

Making Sense of Data

Getting it wrong, and getting it right

John Maindonald

2021-06-17

Contents

| | |
|---|-----------|
| Preface | 7 |
| 1 Overview | 9 |
| 2 Human judgment — two systems | 13 |
| 2.1 The Intuition of Professionals (Kahneman, p. 12) | 14 |
| 2.2 Kahneman overview | 15 |
| 2.3 System 1 and System 2 — further comments | 16 |
| 2.4 Try not to be fooled | 16 |
| 2.5 By classical economic criteria, humans misbehave! | 18 |
| 2.6 Helpful Questions for Life in an Uncertain World | 18 |
| 3 Effective use of graphs | 19 |
| 3.1 General principles | 19 |
| 3.2 Banking — the importance of aspect ratio | 19 |
| 3.3 Varying time intervals — show rates, not counts | 20 |
| 3.4 Helpful web links are: | 21 |
| 4 Selection and survivor bias | 23 |
| 4.1 The hazards of convenience samples | 23 |
| 4.2 UK cotton worker wages in the 1880s | 24 |
| 4.3 How good a guide does the past provide to the future? | 26 |
| 4.4 Tales of standout past success (pp.38–39) | 27 |
| 4.5 Damaged planes — how were the data generated? | 27 |
| 5 Medical and related tests and trials | 29 |
| 5.1 Useful general sources of advice and information | 29 |

| | | |
|-----------|---|-----------|
| 5.2 | Other sources of information | 30 |
| 5.3 | PSA Screening for Prostate Cancer, and more | 31 |
| 5.4 | Randomized Controlled Trials (RCTs) versus Population Studies | 32 |
| 5.5 | Avoid, or expose infants to peanuts? | 34 |
| 5.6 | False Positives (Smith, pp.98-99) | 35 |
| 5.7 | Breast cancer screening — a very contested area | 37 |
| 6 | Fame that may spill over into notoriety | 39 |
| 6.1 | The MMR (measles, mumps, and rubella) scandal | 39 |
| 6.2 | Sally Clark’s disturbing cot death experience | 40 |
| 6.3 | John Snow documents a natural experiment | 41 |
| 6.4 | The Reinhoff and Rogoff saga (Smith, p.55) | 44 |
| | There are further serious issues of interpretation | 47 |
| | Parting comments | 47 |
| 7 | Weighting effects, & the Yule-Simpson Paradox | 49 |
| 7.1 | Deaths from Covid-19 — between country comparisons | 49 |
| 7.2 | The Yule-Simpson paradox — UCB admissions data | 51 |
| 7.3 | Third variables change the story — further examples | 55 |
| 7.4 | Cricket Bowling Averages | 58 |
| 7.5 | Epistatic effects on genetic studies | 58 |
| 8 | Regression and Correlation | 59 |
| 8.1 | What direction does the correlation go? | 59 |
| 8.2 | Regression to the mean | 60 |
| 8.3 | Regression to the mean in a variety of contexts | 63 |
| | Decathlon scores — between event correlations | 63 |
| 8.4 | Moderating subjective assessments | 70 |
| 9 | Covariate adjustments in observational studies | 75 |
| 9.1 | What is driving predictions? — sources of advice | 76 |
| 9.2 | From air, or from water — 1849 deaths from cholera | 77 |
| 9.3 | Are there missing covariates? | 78 |
| 9.4 | The May 2020 Lancet paper that was quickly withdrawn | 79 |
| 9.5 | The uses and traps of “algorithmic” methods – trees | 80 |
| 9.6 | Regression bloopers – examples of other traps | 84 |
| 9.7 | Global mean temperature trends | 87 |
| 10 | The critique of scientific claims | 91 |

| | |
|--|------------|
| <i>CONTENTS</i> | 5 |
| 10.1 What results can be trusted? | 91 |
| 10.2 A look at what can go wrong | 92 |
| 10.3 The case of Eysenck and his collaborators | 93 |
| 10.4 The use of artificial intelligence to detect Covid-19 | 93 |
| 10.5 Laboratory experimental science — what do we find? | 94 |
| 10.6 The Reproducibility: Psychology project | 95 |
| 10.7 When science and commercial interests collide | 96 |
| 10.8 Tricks used to dismiss established scientific results | 99 |
| Further reading — books, videos, and websites | 101 |

Preface

It ain't what you don't know that gets you into trouble. It's what you know for sure that just ain't so. [Advice attributed to Mark Twain]

Data do not stand on their own. It has a context. That context has to be understood and to feed into the way that the data are used, if meaningful conclusions are to be drawn. This booklets aims to identify some of the issues that arise in making sense of data, and to document common misunderstandings.

Using Kahneman's book "Thinking Fast and Slow"¹ as a starting point, we will note mental traps to which humans are prone, and provide examples from the media and from published work.

The hope is that knowledge of the traps will help us to avoid them. Kahneman identifies, as a useful way to think about the issues that arise when humans make judgments, two processes that he calls, respectively, System 1 (jumping) and System 2 (pondering). This is a simplification — but does provide a helpful starting point for discussing the psychology involved.

The questions that data and data analysis are designed to answer can often be stated simply. This may encourage the layperson to believe that the answers are similarly simple. Often, they are not. Be prepared for unexpected subtleties. Always, think carefully how the data were collected, and about how this may limit the conclusions that can be drawn from it.

There is brief attention to issues that arise from practices that have grown up within some areas of laboratory science. Concerns about reproducibility, especially in wet laboratory biology and in psychology, have in the past decade

¹Kahneman (2013)

attracted extensive attention in the pages of *Nature*, *Science*, the *Economist*, psychology journals, and elsewhere. There are self-correcting processes that in the long run work to set the record straight. Better and more consistent attention to independent checks that results can be replicated would work to avoid much of the wasting of effort and time that can occur when later investigators find that the results on which they hoped to build are flawed.

Chapter 1

Overview

Issues and questions that appear repeatedly are:

1. Understand the psychology — why the showcase of fallacies
 - An understanding of the mental traps to which we are prone can help us avoid them.
2. Graphs can reveal surprises. They can also be drawn so that they deceive. Take care!
3. What is the cause? What should we blame, or praise?
 - What causes cholera – bad air or bad water?
 - Does drinking coffee help or harm health?
4. Randomized trials, if done rigorously, are a gold standard. Participants must closely reflect the population to which results will be applied.
5. Population studies, where covariate adjustments are needed to account for prior differences between the two (or more) groups, much more readily mislead?
6. From hindsight to insight — it's easy to be deceived.
 - If a sportsperson is at the top of their form, the only way to go is down. If at the bottom, the only way is up.

- This is true also for success in business.
 - Chance, as well as form, obviously plays a part.
7. The Yule-Simpson paradox appears in many different guises. It arises from what are really weighting effects, because we have naively added numbers in ways that give more weight to some combinations of effects more than to others.
 8. Tall fathers are likely to have tall sons, but shorter than themselves. Tall sons are likely to have tall fathers, but shorter than themselves. This is an example of regression to the mean. What may seem puzzling is that it goes in both directions, from sons to fathers as well as from fathers to sons.
 9. Results become part of established science when they have survived informed critique. Work that claims to revise or overturn established results has itself to survive the same informed critique.
 - There is a great deal of uninformed critique. Watch for the tricks that are used in the attempt to discredit established scientific results that detractors find inconvenient.
 10. Publication processes, widely across many areas of science, urgently require to be updated.
 - They work well where the nature of the work requires input and critique widely from across different disciplinary perspectives.
 - When presented with scientific (or any) claims, ask/check whether claims have been carefully critiqued. For laboratory studies, have results been replicated?

A simple (and, arguably, trick) example

Linda is a 31-year old philosophy graduate, single, outspoken, and bright. As a student, she was deeply concerned with issues of discrimination and social justice, and also participated in anti-nuclear demonstrations. Which of the following is more probable?

- Linda is a bank teller.
- Linda is a bank teller, and active in the feminist movement.

Adding the further descriptor “active in the feminist movement” can only lower

the probability, or just possibly leave it unchanged. Instead of assessing the balance of probabilities, we are tempted to ask which description best meshes with what we have been already told about Linda.

“Linda is active in the feminist movement” is the single descriptor” that respondents see as best fitting Linda. While that was not what was asked, one has to pay close attention to prevent System 1 from substituting that for the question that was asked. Note that the correct answer will be “a bank teller”, irrespective of the way that Linda was characterized before the question was asked.

Chapter 2

Human judgment — two systems

An understanding of the mental traps to which we are prone can help us avoid them. System 1 (jumping) and System 2 (pondering) provide helpful models (one may call them “fictions”), for understanding how humans think.

* System 1 (jumping quickly to a conclusion) responds quickly to immediate dangers, but is prey to priming (suggestive language), framing (response to a choice depends on how it is presented), affect (feeling or emotion), memory illusions, illusions of truth, ... * System 1 may answer an easier question in place of a harder. * System 2 is the process involved when, before making a decision, we think carefully through the issues involved. + Invoke System 2, when time allows, for decisions that matter.

* System 1 and System 2 are both amenable to training * A well-trained System 2 is the key to a better System 1

Further points are

* Untrained humans are poor intuitive statisticians. + Too often, we jump to conclusions, without careful assessment. + Or, we may not be equipped to make an informed and carefully thought through decision. * System 1 and System 2 are both amenable to training + A well-trained System 2 is the key to a better System 1.

[Kahneman: Ch 1; Smith: Ch 1; Nisbett Ch 1-3]

2.1 The Intuition of Professionals (Kahneman, p. 12)

The situation has provided a cue; this cue has given the expert access to information stored in memory, and the information provides the answer. Intuition is nothing more and nothing less than recognition. (Simon, 1992, “What is an Explanation of Behavior?”)

Simon’s comments were in part based on the study of chess masters.

... An expert is one who has built a deep familiarity with the patterns of a given domain and, thus, has a robust body of work from which to reference when making decisions.¹

Obstacles to effective judgment

Effective professional training is designed to ensure that at least some of the results of well-tuned System 2 expert judgment operate at a System 1 level. The biases that Kahneman documents can readily take over, where there should be carefully reasoned System 2 judgment.

- Financial firm executive, explaining why he invested in Ford stock:
 - “Boy, do they know how to make a car!”
 - An easier and related question (do I like Ford cars?) came readily to mind and determined his choice.
- Exemplifies the *affect* heuristic — decisions are guided by like/dislike feelings, with little deliberation or reasoning.

Even those who are experts in their field can be similarly prone to judgments that have no foundation in fact:

“But you’re not an MD like I am. The problem with you is that you do not understand medicine. You see, medical statistics are not like other statistics.”

[Prostate biopsy specialist, in response to evidence that screening is counterproductive. (Levitin, 2015, p.247)]

¹ <http://www.euclidean.com/expert-intuition-and-machine-learning/>

Association of Professional Psychologists, web post on [Arkes and Gaissmaier \(2012\)](#), discussing the response to a U.S. Preventive Services assessment that prostate screening, when used in accordance with then current treatment practices, was doing more harm than good —

Even faced with ... evidence ... [from] a ten-year study of around 250,000 men that showed the test didn't save lives, many activists and medical professionals are clamoring for men to continue receiving their annual PSA test.

New evidence emerges as time proceeds, and there are advances in the approach to treatment. At least part of the problem has been a rush to treatments that themselves risk increasing damage and the risk of death. Note the comment in [Brawley \(2018\)](#) that

Over the past few years, the benefit-to-harm ratio has improved in favor of benefit if the man understands that active surveillance may be a reasonable path if diagnosed.

2.2 Kahneman overview

Part I: Two Systems (System 1 versus System 2)

This aims to show how “... associative memory, the core of System 1, attempts to give a coherent account of what is going on in our world in any instant”

Part 2: Heuristics and biases * Why is it so difficult for us to think statistically?

* We easily think associatively, ... metaphorically, casually, but statistics requires thinking about many things at once ...

Part 3: Choices * We have “an excessive confidence in what we think we know”

Part 4: A conversation with the discipline of economics * The Prospect theory model of choice + “... unfortunate tendency to treat problems in isolation” + Framing effects — inconsequential features shape choices + A challenge to the assumptions of standard economics.

Part 5: Two selves — Experiencing, remembering * These do not always have the same interests. * Automatic memory formation has its own rules + Exploit to improve the memory of a bad episode * The virtues of educating gossip +

Use to improve the quality of judgments and decisions + Judge decisions by how they are made, not by outcome.

2.3 System 1 and System 2 — further comments

- System 1 Features
 - It may substitute an easier question
 - It responds to priming, framing, affect, memory illusions, illusions of truth, ...
 - It has little understanding of logic and statistics
 - It cannot be turned off
 - When flummoxed, it calls on System 2
- The world of System 2
 - It keeps you polite when angry, alert when driving at night
 - In its world, gorillas do not cross basket-ball courts ...
 - There are some vital tasks that only System 2 can perform
 - Problems that put 1 & 2 in conflict may require large mental effort & self-control to overcome the impulses and intuitions of System 1
- Healthy living is a compromise
 - Recognize situations where mistakes are likely
 - Aim to avoid significant mistakes when the stakes are high

Kahneman describes System 1 and System 2 as “fictions”! I prefer to describe them as providing a model that is helpful in accounting for the way that humans make judgments.

- Human mental quirks make it useful to think of these as agents.
- Stories about active agents who have personalities, habits, and abilities readily take root in our minds
- Ascribing calculation to System 2 is shorthand for ...
 - “This requires mental effort. Do not do when driving!”
 - Pupils will be dilated, and heart rate elevated.

2.4 Try not to be fooled

- We make judgments based on evidence that is too limited

- We are easily fooled by irrelevancies
 - Kahneman has brought together evidence on what & how.
- Getting & assessing good evidence takes discipline & work
 - Data may (rarely) be there for the taking; but someone has to notice, and collate the data.
 - Randomized experiments are often the ideal, but require meticulous planning. If differences are small, the numbers required may be very large. See further, Subsection 5.4.
- Keep in mind Yule-Simpson “paradox”, which we will encounter later in Subsection 7
 - The paradox lies in the failure of human intuition to accommodate straightforward arithmetic!
 - <https://www.youtube.com/watch?v=ZDinnCwP3dg> (This animated video provides a short overview)

An irony is that Kahneman was, as he has acknowledged, himself fooled into taking at face values papers that claimed to show that verbal concepts could have the effect of altering behaviour. Thus

- Being asked to write down stories about unethical deeds made people more likely to want to buy soap;
- Subtly drawing attention to money, e.g., leave banknotes lying around, made people feel more self-sufficient, and care less about others;
- Priming people with old age related words leads people to walk more slowly away from the lab as research assistants armed with stopwatches timed their movements.

As [Ritchie \(2020\)](#) notes (p.28) Kahneman was not alone in being fooled — the study about priming with old age related words has been extensively cited in psychology textbooks. None of these claims have stood up in attempts at replication, with larger numbers and with greater care to avoid unconscious sources of bias. Thus, in the replication of the study relating to age-related words, infrared beams were used to measure time taken to walk between two points in a hallway, rather than research assistants who knew the group to which participants had been assigned.

Think again — a very simple example

Is symptom X associated with disease A?

| | Has Disease | No Disease |
|-------------------|-------------|------------|
| Symptom X Present | 20 | 10 |
| Symptom X Absent | 80 | 40 |

The symptom occurs with the same relative frequency, whether or not a person has the disease. Nisbett comments that most people, including nurses and doctors, interpret such evidence wrongly (Nisbett, 2016, pp.129-130).

2.5 By classical economic criteria, humans misbehave!

The discipline of “behavioural economics” largely took shape as a result of the work of Richard Thaler. Kahneman was one of two mentors who strongly influenced Thaler — the other was Amos Tversky, who had worked closely with Kahneman.

Thaler and Ganser (2015) explores the extent to which humans do not behave like the rational agents of classical economics, agents to whom Thaler gives the name “econs”. Added to the irrationality with which we often act is that our personal priorities are unlikely to align precisely with those of econs.

2.6 Helpful Questions for Life in an Uncertain World²

1. Risk of what? (Showing a symptom, ..., Death)
2. What is the time frame? (next 10 years, or lifetime)
3. How big is the risk? (Look at risk in absolute terms)
4. Does the risk apply to me? (Age, sex, health, ...)
5. What are the harms of “finding out”? (False alarms, invasive diagnostic procedures, unnecessary or dangerous treatments.)

²These came originally from the Harding Center web site <https://www.hardingcenter.de/en>.

Chapter 3

Effective use of graphs

3.1 General principles

- Focus the eye on features that are important
- Avoid distracting features
- Lines that are intended to attract attention can be thickened
- Where points should be the focus, make them large & dark
 - It often makes sense to de-emphasize the axes.
- If points are numerous and there is substantial overlap, use open symbols, and/or use symbols that have some degree of transparency.
- Different choices of color palettes are effective for different purposes.
- Bear in mind that the eye has difficulty in focusing simultaneously on widely separated colors that are close together on the same graph.

3.2 Banking — the importance of aspect ratio

Patterns of change on the horizontal scale that it is important to identify should bank at an angle of roughly 45° above or below the horizontal. Angles in the approximate range 20° to 70° may be satisfactory, and the aspect ratio should be chosen accordingly.

```
par(mgp=c(2,.5,0))  
plot((1:50)*0.92, sin((1:50)*0.92), xlab="", ylab="", fg='gray')
```

```

mtext(side=3, at=0.75, "A: Pattern is hidden to view", adj=0, cex=1.25)
plot((1:50)*0.92, sin((1:50)*0.92), xlab="", ylab="", asp=2)
mtext(side=3, at=0.75, "B: Pattern is now obvious", adj=0, cex=1.25)

```

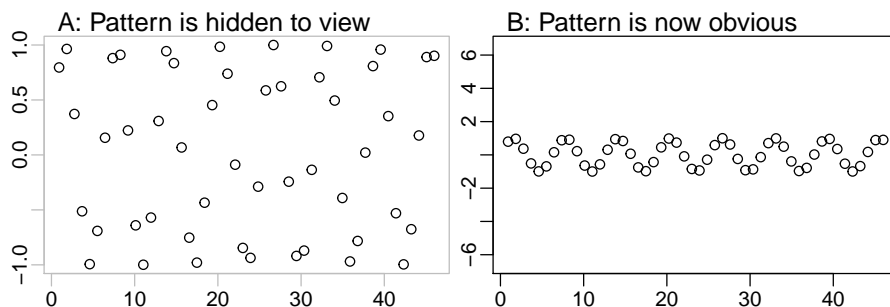


Figure 3.1: The same data are used for both graphs. The pattern that is not obvious in Panel A is very obvious in Panel B

3.3 Varying time intervals — show rates, not counts

```

par(mar=c(4.1,3.1,2.1,3.1),mgp=c(2.25,0.5,0))
n2 <- c(1,1,4,12,15,29,27,11)
n1 <- c(1,1,4,12,15,29,27,39)
plot(1:8, n2, xaxt="n", xlab="", ylim=range(n1),
     ylab="Number of Nobel prizes", type="l", fg='gray')
axis(1, at=1:8, labels=paste0(seq(from=1901, to=1971, by=10), "-"), cex.axis=0.9)
axis(1, at=1:8, labels=paste0(c(seq(from=1910, to=1970, by=10), 1974)),
     line=0.8, lty=0, cex.axis=0.9)
axis(1, at=1:7, labels=FALSE)
mtext(side=3, line=0.5, at=0.5, "Total prizes (black line)", adj=0, cex=1.25)
mtext(side=3, line=0.5, at=8.5, "Average prizes per year (blue dots)", col='blue')
points(1:8, n1, lwd=8, col=adjustcolor('blue', alpha=0.5), pch=1)
lines(7:8, tail(n1,2), lty=2,lwd=2, col=adjustcolor('blue', alpha=0.5))
axis(4, at=(0:4)*10, labels=paste(0:4), col.ticks=adjustcolor('blue', alpha=0.5))
mtext(side=4, line=1.8, "Nobel prizes per year", col='blue')

```

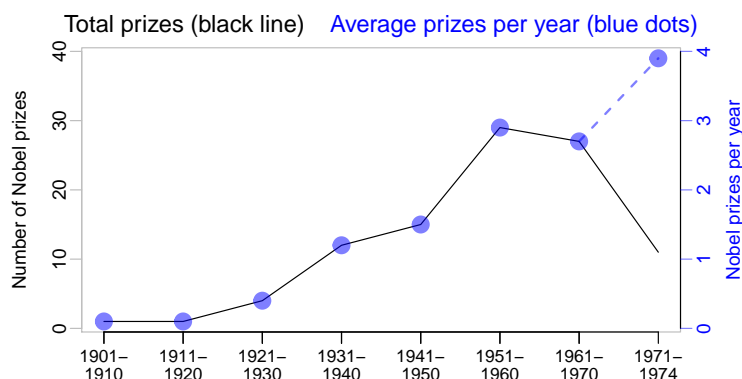


Figure 3.2: The black line shows numbers of US Nobel prizes, for given time intervals. The blue dots show average number per year, as indicated by the axis on the right.

A graph that was essentially the black segmented line in Figure 3.2 appeared in [National Science Foundation \(1975\)](#) “Science Indicators, 1974”. The segmented line gives a highly misleading impression for the four years 1971-1974, as opposed to earlier points, where numbers are totals over decades. It joins up a final point that is a different measure from earlier points.

A proper comparison between the number of Nobel prizes in the earlier decades and the number in the four-year period 1971-1974 adjusts for the length of time (ten years versus four years) to which the count relates. The blue dots, and the axis on the right, fairly represent change over time.

The same principle applies for intervals of measures other than time — for example of length or volume.

3.4 Helpful web links are:

- Good & bad graphs (Murrell, lecture notes)¹
- A short tour of bad graphs (Schwartz)²
- Misleading graphs³

¹<https://www.stat.auckland.ac.nz/~ihaka/120/Lectures/lecture03.pdf>

²www.stat.sfu.ca/~cschwarz/Stat-650/Notes/PDF/ChapterBadgraphs.pdf

³www.statisticshowto.com/misleading-graphs/

- Color choices⁴
- Color Brewer⁵

⁴[<www.stonesc.com/pubs/Expert%20Color%20Choices.pdf>](http://www.stonesc.com/pubs/Expert%20Color%20Choices.pdf)

⁵[<www.statisticshowto.com/misleading-graphs/>](http://www.statisticshowto.com/misleading-graphs/)

Chapter 4

Selection and survivor bias

In mind here are cases where the data are not a random sample.

4.1 The hazards of convenience samples

Quota sampling has often been used as an alternative to random sampling — quotas are set for age categories, male/female, and socioeconomic categories that are designed to ensure that the sample is representative of the wider population. In polls prior to the 1948 US presidential election that pitted democrat Harry Truman against republican Thomas Dewey, pollsters were given strict quotas, but otherwise left free to decide who they would approach. Polls by three different organizations gave Dewey a lead of between 5% and 15%. In the event, Truman led by 5%.

4.1.1 Convenience samples sometimes have a story to tell

This is no to rule out all use of convenience samples. Convenience samples, taken within a limited population, can sometimes be useful in setting a lower bound. It is strongly in the public interest that scientists have reasonable freedom for responsible expression of their minds on issues of public concern. In an informal 2015 survey, 151 Crown Research Institute scientists (out of 384 who responded)

answered yes to the question “Have you ever been prevented from making a public comment on a controversial issue by your management’s policy, or by fear of losing research funding?” The 384 who responded will undoubtedly be a biased sample. Irrespective of the size of the bias, the number who had not been allowed to speak their mind was large enough to be a cause for serious concern. Hon Joyce’s response, to the effect that as this was not a scientific survey of all CRI scientists (to this extent, true), its evidence of large concern could be ignored, was an evasion. Equally disturbing was the reaction of the NIWA management, suggesting that they did not accept a responsibility to defend transparency.¹

4.2 UK cotton worker wages in the 1880s

Prior to the [Boot and Maindonald \(2008\)](#) paper² the main source of published information on cotton worker wages in the UK in the late 19th century were results from an 1889 US Bureau of Labor survey, intended for use for comparison with the US cotton industry wages. Figure 4.1 compares the US Bureau of Labor survey numbers with 1886 census numbers of different types of full time UK cotton operatives.

```
cap22 <- "Cotton worker wages in the UK --- 1889 US Bureau of Labor 'survey'
versus 1886 census data. Wages are in pence per week."
```

```
## detach("package:latticeExtra")
suppressPackageStartupMessages(library(ggplot2, quietly=TRUE))
cottonworkers <- DAAG::cottonworkers
names(cottonworkers)[3] <- "Av_wage"
cottonworkers$occ <- row.names(cottonworkers)
cottonworkers$occ[cottonworkers$census1886<2500] <- ""
cottonworkers$occ2 <- row.names(cottonworkers)
cottonworkers$occ2[!cottonworkers$census1886 %in% c(208, 1250:2500)] <- ""
gph <- ggplot(cottonworkers,
              aes(census1886,survey1889,label=occ,
                  size=Av_wage, hjust=1))+
```

¹See <https://sciblogs.co.nz/infectious-thoughts/2015/08/28/niwa-in-astonishing-attack-on-scientist-association/>

²“New estimates of age- and sex- specific earnings and the male-female earnings gap in the British cotton industry, 1833-1906”


```

geom_abline(intercept=0, slope=0.00986, color="gray") +
geom_point()+
geom_text(size=4,nudge_x=-180) +
geom_text(aes(y=survey1889+rep(c(5,-5),c(7,7)), label=occ2),
          size=4, hjust=0, nudge_x=100) +
  xlab("Number in 1886 census") + ylab("Number in 1889 survey")
leg <- matrix(c("The line shows what 1899",
                "survey numbers would be,",
                "if relative numbers in",
                "worker categories were",
                "as in the 1886 census"), ncol=1)

library(gridExtra)
tt <- ttheme_minimal(base_size=12,
                    core = list(fg_params=list(hjust = 0, x=0)))
gph + annotation_custom(tableGrob(leg, theme=tt),
                        xmax=3600, ymin=100, ymax=150)

```

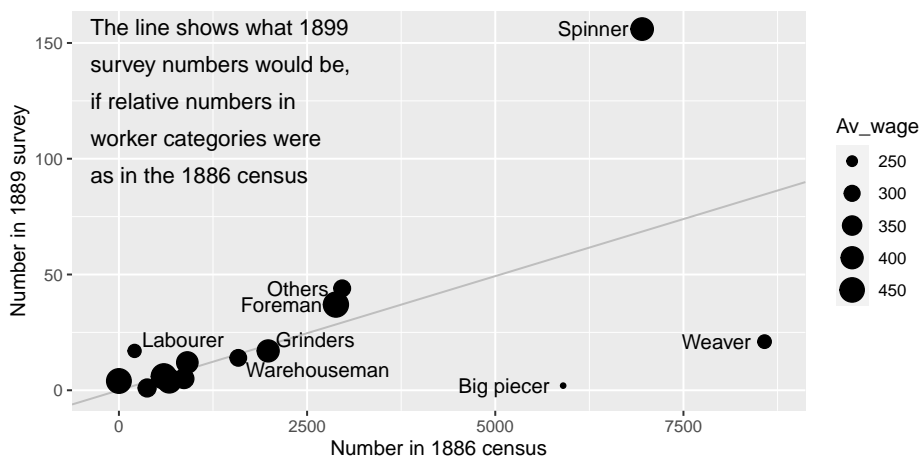


Figure 4.1: Cotton worker wages in the UK — 1889 US Bureau of Labor ‘survey’ versus 1886 census data. Wages are in pence per week.

The 1889 survey shows some strong biases — a result it would seem of geographical bias and of the informal data collection methods that were used. The high

wages given to spinners were grossly over-weighted in the US Bureau of Labor survey, while Big Piecers and Weavers were grossly under-represented. A guess is that workers were asked for information on their wages as they left work, and that the survey personnel happened to catch employees at a time when there was a large preponderance of highly paid spinners, and an untypically small number of big piecers and weavers. The net effect was a gross over-estimate of average wages.

4.3 How good a guide does the past provide to the future?

There is a mix of selection and survivor bias arise when data from the past are used as a guide to the future, with no allowance for the source/target difference. The target about which we wish to make judgments lies in the future, while the data are from the past. Think about a business that is planning for the future. One can never know, until after the event, all the ramifications of the choices made.

Businesses may be selected as examples of effective business practice because they were, at the time when the data were collated, successful. Likewise, it is athletes who have been successful in the recent past who are likely to be selected to appear on the covers of sports magazines. In either case, this gives a biased picture of what can be expected in the future — in many cases those who are picked out will be close to the peak of their success, and/or have had unusual luck. High levels of success in the recent past will not always translate into success in the following year or years. How often will past success translate into future success? In order to discover, it is necessary to collect relevant data.

Subsection 4.5 will discuss a widely cited example of survivor bias that relates to the relative density of bullet holes, in different parts of the plane, in planes that had returned from battle in World War II. The key to interpreting the data was to recognize that these were the planes that had survived.

4.4 Tales of standout past success (pp.38–39)

Collins (2001) — Good to Great

- From 1,435 companies, Collins identified 11 as standouts
 - Since 2001, 5 better than average; 6 worse
 - Fannie Mae — 2001: >\$80 per share: 2008: <\$1
 - Circuit City: Bankrupt in 2009

Waterman and Peters (1982) — In Search of Excellence

- 43 successful companies:
 - “Bias for action”; “Close to consumer”
- From the 35 that were publicly listed:
 - Since 1982: 15 better than average; 20 worse

In all cases, companies that were chosen as examples of standout success were likely to be near to the peak of their performance, as judged by Collins, or by Waterman and Peters. Overlaid on this is the regression effect that will be discussed in Chapter 8.

4.5 Damaged planes — how were the data generated?

Abraham Wald’s insight was that survivor bias was to be expected, with the density of bullet holes providing evidence about the extent of bias, and the implications for identifying the part(s) of the planes that would benefit most from additional protection. See [Abraham Wald and the Missing Bullet Holes](#)³, which is an excerpt from [Ellenberg \(2015\)](#).

The numbers of hits per square foot were:

```
c('Engines'=1.11, "Fusilage"=1.73, "Fuel system"=1.55, "Rest of plane"=1.8)
```

³<https://medium.com/@penguinpress/an-excerpt-from-how-not-to-be-wrong-by-jordan-ellenberg-664e708cfc3d>

| ## | Engines | Fusilage | Fuel system | Rest of plane |
|----|---------|----------|-------------|---------------|
| ## | 1.11 | 1.73 | 1.55 | 1.80 |

It was not possible to fit protective plates everywhere, as this would have made the planes overly heavy? Wald argued that the holes very likely were spread pretty much uniformly over the planes as a whole, those that were shot down, and those that survived, The reason for relatively fewer bullet holes in the engine and fuel system areas was that hits in those areas were more likely to bring the plane down, so that they did not return.

Chapter 5

Medical and related tests and trials

5.1 Useful general sources of advice and information

The Harding Center for Risk Literacy

There is extensive informative content on the [Harding Center for Risk Literacy site](https://www.hardingcenter.de/en)¹ site. Note, in particular:

- [Team](https://www.hardingcenter.de/en/team)²
- [Medical Fact Boxes](https://www.hardingcenter.de/en/medical-fact-boxes)
- [Publications](https://www.hardingcenter.de/en/publications)
- [In the Media](https://www.hardingcenter.de/en/in-the-media)

¹<https://www.hardingcenter.de/en>

²<https://www.hardingcenter.de/en/team>

Harding Center Medical Fact Boxes

Medical fact boxes³ provide visual and tabular summaries of the current “best” evidence, from randomized controlled trials. The comparison may be with a placebo, or with an alternative that is known to be effective. Detailed references are given. Where available, reliance is on Cochrane studies. Some of the available fact boxes are:

- Early cancer detection: breast, ovarian, prostate, colon, cervix
- Cardiovascular Diseases
- Vaccines
- Low back pain
- Osteoarthritis of the knee
- Dietary Supplements - Selenium

Note also

- information material on SARS-CoV-2 and Covid-19⁴
- Are you risk literate? — Risk quiz⁵

5.2 Other sources of information

- Winton Centre for Risk and Evidence Communication⁶
- [US Consumer health information site⁷

The Cochrane Center⁸

The Cochrane Center’s mission is “to promote evidence-informed health decision-making by producing high-quality, relevant, accessible systematic reviews and other synthesized research evidence.” They rely heavily on meta-analyses, looking for the balance of evidence across all relevant studies.

³<https://www.hardingcenter.de/en/projects-and-collaborations/fact-boxes>

⁴<https://www.hardingcenter.de/en/transfer-and-impact/what-you-should-know-about-sars-cov-2-and-covid-19>

⁵https://docs.google.com/forms/d/e/1FAIpQLSfw3Pe5a_g6lALY7ad3Aq7fBu2Sro4uteUKqvvyvCARICxOQ9g/viewform?usp=sf_link

⁶<https://wintoncentre.maths.cam.ac.uk/>

⁷<https://medlineplus.gov/>

⁸<https://www.cochrane.org/>

5.3 PSA Screening for Prostate Cancer, and more

Numbers (rounded) in the following table are from a Harding Center fact box. They are for men 50 years or older who either did or did not participate in prostate cancer screening, using the PSA test, for 16 years.⁹

| | 1000 men, No screening | 1000 men, Screening |
|--------------------------|------------------------|---------------------|
| Deaths (any cause) | 322 | 322 |
| Deaths (prostate cancer) | 12 | 10 |
| Biopsy & false alarm | 0 | 155 |
| Unnecessary treatment | 0 | 51 |

About 10 out of every 1,000 men with screening, and 12 out of every 1,000 men without screening died from prostate cancer within 16 years. This means that 2 out of every 1,000 people could be saved from death from prostate cancer by early detection using PSA testing. This was not reflected in overall mortality.

Numbers for benefits are based on four studies with about 77,000 participants (progressive cancer), four studies with about 472,000 participants (overall mortality), and eleven studies with about 619,000 participants (prostate cancer specific mortality). The numbers for harms are based on seven studies with approximately 128,000 participants (false-positive results within three to six participations in PSA testing for early detection) and nine studies with approximately 274,000 participants (over-diagnosis and over-treatment). See the web site for references to the studies.

Unlike the biopsies that may follow a positive PSA test, the PSA test has no direct potential to cause physical harm. The harm results from an undue readiness to use the test result as a reason for further testing and treatment that will in many cases itself cause harm. “Wait and watch” appears the preferred strategy.

See also: Levitin (2015), Chapter 6, and [How Patients Think, and How They](#)

⁹<https://www.hardingcenter.de/en/early-detection-of-cancer/early-detection-of-prostate-cancer-with-psa-testing>

Should¹⁰

5.4 Randomized Controlled Trials (RCTs) versus Population Studies

Two types of study are widely used in medical and other contexts — randomized controlled trials, and population-based studies. These can, in both cases, be broken down into further sub-types. There may be elements of both these types of studies.

Randomized Controlled Trials (RCTs) — the gold standard?

- Ensure that apples are compared with apples
 - Use a random mechanism to assign to treatment as against control — in a medical screening study to “screen” or “not screen”
 - Treatment and control must otherwise be treated in the same way.
- There must be strict adherence to a protocol
 - Minor departures that may, e.g., allow unconscious bias in the way that results from the different groups of participants are measured, can invalidate results.
- Results apply, strictly, to those who meet the trial entry criteria
 - This may limit relevance to general population

Especially in medical trials, think carefully about the outcome measure.

- In a screening trial, e.g., for prostate cancer, there are risks both for those who test positive, and for those who test negative.
 - The process used to check for cancer may itself bring a smaller or larger element of risk.
 - Positives may be false positives, leading to more invasive checks which may themselves carry a risk. Thus, for prostate cancer, a positive PSA test is likely to lead to a biopsy that itself has been es-

¹⁰<https://www.nytimes.com/2011/10/09/books/review/your-medical-mind-by-jerome-groopman-and-pamela-hartzband-book-review.html>

5.4. RANDOMIZED CONTROLLED TRIALS (RCTS) VERSUS POPULATION STUDIES33

- Some slow growing cancers may be better left untreated, rather than exposing the patient to a treatment that may itself do serious damage.
- Thus, think carefully about the choice of outcome measure. It is not enough to show that a screening program will pick up otherwise undiagnosed cancers.

For a helpful animated summary of some of the key issues, see:

<https://www.youtube.com/watch?v=Wy7qpJeozeC>

A note in passing: HiPPO decisions vs A/B testing

Randomized studies are widely used outside of medicine. Randomization is a key component of the way that Google and others test out, e.g., the effect of different web page layouts.

- HiPPO = “Highest paid person in the Office.”
- The term “A/B testing” is sometimes used to refer to randomized testing of alternatives.

A/B testing helped propel Obama into office! An experiment was conducted that involved 15 million people, or about 25%, from its email list. The signup forms had one of nine different combinations of images with words on which recipients were invited to click, thus:

| | Learn more | Join us | Sign up now |
|--------------------------|------------|---------|-------------|
| Obama photo | ✗ | ✗ | ✗ |
| BW photo of Obama family | ✓ | ✗ | ✗ |
| Obama speaking | ✗ | ✗ | ✗ |

The black and white photo of the Obama family, with the words “Learn more”, generated the most clicks.

Young (2014) gives an account of A/B testing as it might be used for improving library user experience.

Population studies — groups must be broadly comparable

- Adjust prunes to look like apples (is it possible?)

- Can one ever be sure that the adjustments do the job?
- Potential for biases is greater than for RCTs

Where a treatment is compared with a control group, the idea is to use a regression type approach to adjust for differences in such variables or factors as age, sex, socioeconomic status, and co-morbidities. “Propensity score” approaches try to summarize such group differences in a single variable (or, in principle, two or more variables) that measure the propensity to belong to the treatment as opposed to the control group. While their effectiveness for this purpose may be doubted, they can be used to provide insightful graphs that check the extent to which the groups are broadly comparable on the variables and/or factors used to adjust for differences.

Issues for all types of study

What are the relevant outcome measures?

- e.g., cancer – malignancies found & removed, or deaths
 - deaths from cancer, or from all causes (for some individuals, the treatment may be more damaging in its medium to long term effect than the cancer)

Care is required to deal with survivor, as well as other, biases.

5.5 Avoid, or expose infants to peanuts?

Clinical practice guidelines introduced in or around the year 2000 had “recommended the exclusion of allergenic foods from the diets of infants at high risk for allergy, and from the diets of their mothers during pregnancy and lactation.”

It was then a surprise to find that the prevalence of peanut allergy has substantially increased in the recent past, doubling in Europe between 2005 and 2015, suggesting that advice given to parents of young children to avoid foods containing peanuts may have been counterproductive. This reassessment was supported, at least for infants who at four months had either severe eczema or food allergy or both, and thus were at high risk of developing a peanut allergy, by the LEAPS study reported in [Du Toit et al. \(2015\)](#).

As noted, the LEAPS study was limited to infants who at four months had either severe eczema or food allergy or both.

Infants were stratified into two groups following a skin-prick test, with each group then randomized between those exposed to peanut extract, and those not exposed.

Among 530 infants in the population who initially had negative results on the skin-prick test, the prevalence of peanut allergy at 60 months of age was 13.7% (37/270) in the avoidance group and 1.8% (5/272) in the consumption group.¹¹ Among the 98 participants who initially had positive test results, the prevalence of peanut allergy was 35.3% (18/51) in the avoidance group and 10.6% (5/47) in the consumption group. There was no between-group difference of consequence in the incidence of serious adverse events.

In both groups, numbers and percentages are for those who were assigned to the group and whose results could be evaluated, whether or not they followed the treatment protocol to which they were assigned. In technical terms, these are results from an “intention to treat” analysis. Such an analysis is designed to mirror what can be expected in practice — not everyone who starts off in one group will stick to it. It answers questions about what to do with subjects who did not fully follow the treatment to which they were assigned.

The results were followed, in 2016, by changes to guidelines that recommended introduction of peanut and other allergenic foods before 12 months. The assumption that avoiding early exposure to peanuts would reduce risk of later development of peanut allergy was, it was judged, likely wrong for all infants.

5.6 False Positives (Smith, pp.98-99)

In contexts where the number of false positives is likely to be high relative to the number of true positives, screening programs may have serious downsides that outweigh the benefits.

- A test for marijuana has a 95% accuracy (true positive) rate
- Of those who test positive, what fraction are marijuana users?

Administer test to 10,000 employees — 500 use, 9500 do not

¹¹There were twelve further infants in this group whose results could not be evaluated.

```
tab <- matrix(c(475,475,950,25,9025,9050,500,9500,10000), ncol=3, nrow=3)
dimnames(tab) <- list(c("User", "Not user", "Total"), c("Test +ve", "Test -ve", "Total"))
tab <- xtable::xtable(tab, digits=0)
print(tab, type='latex', comment=FALSE, floating=FALSE, scalebox=1.0)
```

| | Test +ve | Test -ve | Total |
|----------|----------|----------|-------|
| User | 475 | 25 | 500 |
| Not user | 475 | 9025 | 9500 |
| Total | 950 | 9050 | 10000 |

- Consider a test for excess iron syndrome that has ~80% accuracy
Excess iron syndrome has the name “haemochromatosis”. Again, administer it to 10,000 individuals, where 50 have the syndrome, 9950 do not (a rate of 1 in 200).

```
tab <- matrix(c(40,1990,2030,10,7960,7970,50,9950,7000), ncol=3, nrow=3)
dimnames(tab) <- list(c("Haem...", "Not haem...", "Total"), c("Test +ve", "Test -ve", "Total"))
tab <- xtable::xtable(tab, digits=0)
print(tab, type='latex', comment=FALSE, floating=FALSE, scalebox=1.0)
```

| | Test +ve | Test -ve | Total |
|-------------|----------|----------|-------|
| Haem... | 40 | 10 | 50 |
| Not haem... | 1990 | 7960 | 9950 |
| Total | 2030 | 7970 | 7000 |

The Aspirin story (randomized trials)

22,000 males, 1st 5 years of study: <http://nyti.ms/2rgjFO2>

| | Placebo group | Aspirin group |
|-------------------------|---------------|---------------|
| Fatal heart attacks | 18 | 5 |
| Non-fatal heart attacks | 171 | 99 |

Later work: Benefits are limited to those with previous symptoms

2009 study: <https://www.theguardian.com/society/2009/aug/31/aspirin-britis-h-heart-foundation>

Major bleeding episodes, requiring admission to hospital:
were 20 (1.2%) in the placebo group; 34 (2%) in aspirin group

Summary of current evidence: <http://mayocl.in/1cBM0ze>

5.7 Breast cancer screening — a very contested area

See, e.g., [Raichand et al. \(2017\)](#). The review starts with the comment:

The recent controversy about using mammography to screen for breast cancer based on randomized controlled trials over 3 decades in Western countries has not only eclipsed the paradigm of evidence-based medicine, but also puts health decision-makers in countries where breast cancer screening is still being considered in a dilemma to adopt or abandon such a well-established screening modality.

As with PSA screening for prostate cancer, the [Harding Center Fact Box for Mammography Screening](#)¹² is a good place to go for an up to date summary of the evidence. Those who are convinced of the virtues of screening may wish to challenge the evidence as presented at the time of writing. On what grounds, however? It may be argued that

- Improvements in medical procedures since the time of the trials on which the Fact Boxes are based may have changed the balance of risk.
- The trials do not apply to New Zealand conditions.

Both cases require justification.

¹²<https://www.hardingcenter.de/en/early-detection-breast-cancer-mammography-screening>

Chapter 6

Fame that may spill over into notoriety

6.1 The MMR (measles, mumps, and rubella) scandal

Andrew Wakefield was the lead author of a study published in 1998, based on just twelve children, that claimed to find indications of a link between the MMR (measles, mumps, and rubella) vaccine and autism. The journalist Brian Deer had a key role in identifying issues with the work, including fraudulent manipulation of the medical evidence.

- Wakefield had multiple undeclared conflicts of interest
- Funding came from a group of lawyers who were interested in possible personal injury lawsuits
- It emerged that from 9 children said to have regressive autism
 - Only 1 had been diagnosed; 3 had no autism
 - 5 had developmental problems before the vaccine

Wakefield's 1998 claims were widely reported

- Vaccination rates in the UK and Ireland dropped sharply

- The incidence of measles and mumps increased, resulting in deaths and in severe and permanent injuries.

Wakefield was found guilty by the General Medical Council of serious professional misconduct in May 2010 and was struck off the Medical Register.

Following the initial claims in 1998, multiple large epidemiological studies failed to find any link between MMR and autism.

Fact boxes on the Harding site summarize evidence of the effectiveness of the [MMR vaccine](#)¹

6.2 Sally Clark’s disturbing cot death experience

Sudden Infant Death (SID), also referred to as “cot death”, is the name given to the unexplained sudden death of very young children. The story of Sally Clark’s unfortunate brush with the law, following the death of a second child, is interesting, disturbing, and educational. Sally’s experience highlighted ways in which the UK legal system needed to take on board issues that affect the use of medical evidence — issues of a type that can be important in medical research.

Meadow argued that “unless proved otherwise, one cot death is tragic, two is suspicious and three is murder.” Was this an example of “Too hard! Try something easier, and wrong!” Or was it the triumph of assumed “knowledge” over hard evidence?

Following the “cot death” of the second of Sally Clark’s two children, Meadow gave evidence at her trial in 1999, and appeal in 2000. Sally Clark was finally acquitted in 2003.

- Meadow gave 1 in 8,500 as cot death rate in affluent non-smoking families
 - Squared 8,500 to get odds of 73,000,000:1 against for two deaths.
 - Meadow assumed, wrongly, that the probability of a death from natural causes was the same in all families.
 - A first “cot death” is, in some families at least, evidence of a greater proneness to death from natural causes.
- Royal Statistical Society press release: “Figure has no statistical basis”

¹<https://www.hardingcenter.de/de/search/node?keys=mmr>>

- 2000 appeal judges: The figure was a “sideshow” that would not have influenced the jury’s decision.
- The appeal judges’ statement was described by a leading QC, not involved in the case as: “a breathtakingly intellectually dishonest judgment.”
- 2003: It emerged that 2nd death was from bacterial infection
 - In a second appeal, Sally Clark was freed.
- 2005: Meadow was struck from medical register
 - 2006: reinstated by appeal court — misconduct fell short of “serious”!

Meadow was in effect assuming, without evidence, the absence of family-specific genetic or social factors that make cot deaths more likely in some families than in others. Why was Meadow’s assumption of independence not challenged in the 1999 trial? Issues of whether or not different pieces of evidence are independent are surely crucial to assessing the total weight of evidence. Meadow had only enough understanding of probability to be dangerous.

Sally Clark was finally acquitted on her second appeal in 2003, with her sense of well-being damaged beyond repair. The forensic evidence had been weak. The web page <https://plus.maths.org/content/os/issue21/features/clark/index> has a helpful summary of the statistical issues. It quotes a study that suggests that the probability of a second cot death in the same family is somewhere between 1 in 60 and 1 in 130. Even after this adjustment, the probability of death from natural causes of two children in the same family is low. But so also is the probability that an apparently caring mother from an affluent middle class family, with no history of abuse, will murder two of her own children. Those are the two probabilities that must be compared. Anyone who plans to work as a criminal lawyer ought to understand this crucial point. In a large population, there will from time to time be two deaths from natural causes in the one family.

6.3 John Snow documents a natural experiment

The history that led up o the 1854 “experiment”.

- 6,500 died from cholera in London in 1832
 - Medical opinion blamed “miasma”, or noxious air
(Stink from rotting garbage, faeces, & pollution in Thames)

- Poor areas had more cholera – worst smell, sanitation
(But, also, people were older, houses poorly heated, ...)
- Cesspits – night-soil periodically taken away
- 1842: [Edwin Chadwick, in The Sanitary Conditions of the Labouring Population \(1842\)](#)² showed a direct link between poor living conditions, disease and life expectancy
 - Like others who accepted the “miasma” theory of disease, Chadwick did not understand that cholera, along with some other major diseases, were water-borne.
- 1848: the Nuisances Removal and Diseases Prevention Act (Gazette issue 20637) was passed with the aim of stopping the 1848-9 epidemic. The suggestion was to effectively dump the contents of cesspools and raw sewage pits into the Thames, which was London’s main source of drinking water. This only served to exacerbate the problem.
- 1848-49 epidemic followed shortly after the cesspits were banned
 - Both water companies — Lambeth, and Southwark and Vauxhall, were taking water from the same polluted source.
 - Death rates were high for both companies
- 1850: Arthur Hassall’s careful microbiological study:
 - “... a portion of the inhabitants are made to consume
... a portion of their own excrement, and ... to pay
for the privilege” (Hassall, 1850)

See [Cholera epidemics in Victorian London](#)³

The 1854 epidemic — a natural experiment

The 1852 act required water supply companies to move water intake upriver by 1855. By the time of the 1854 epidemic, Lambeth had moved the intake 22 miles upriver, where the water was not contaminated by London sewage. The Southwark and Vauxhall intake was unchanged until 1855. Data on the distribution of cholera in the 1854 epidemic then allowed Snow to test the claims made in his 1849 study.

“The experiment, too, was on the grandest scale. No fewer than 300,000 people ..., from gentles down to the very poor, were divided into two groups without

²<https://www.sciencemuseum.org.uk/objects-and-stories/medicine/cholera-victorian-london>

³<https://www.thegazette.co.uk/all-notices/content/100519>

their choice, and, in most cases, without their knowledge; one group being supplied with water containing the sewage of London, and, amongst it, whatever might have come from the cholera patients, the other group having water quite free from such impurity.”

Lambeth versus Southwark & Vauxhall

```
tab <- cbind("#Houses"=c(40046,26107,256423),
             "#Deaths"=c(1263,98,1422),
             "Rate per 10,000"=c(315,37,59))
dimnames(tab)[[1]] <- c("Southwark & Vauxhall","Lambeth","Rest of London")
print(xtable::xtable(tab, auto=TRUE),type='latex', comment=FALSE, floating=FALSE, s
```

| | #Houses | #Deaths | Rate per 10,000 |
|----------------------|---------|---------|-----------------|
| Southwark & Vauxhall | 40046 | 1263 | 315 |
| Lambeth | 26107 | 98 | 37 |
| Rest of London | 256423 | 1422 | 59 |

Use water from the brewery, and stay healthy!

Snow noted that “Within 250 yards of the spot where Cambridge Street joins Broad Street there were upwards of 500 fatal attacks of cholera in 10 days...”. By contrast, none of the employees of a local Soho brewery developed cholera.

The reason, he judged, was that they drank water from the brewery (which had a different source from the Broad St pump) or just drank beer alone.

[Coleman \(2019\)](#) gives detailed comments on Snow’s work. It took a further ten years for the medical establishment to begin to accept Snow’s conclusions.

```
library(HistData)
SnowMap(polygons=TRUE, main="Snow's Cholera Map with Pump Polygons")
```

Snow's Cholera Map with Pump Polygons

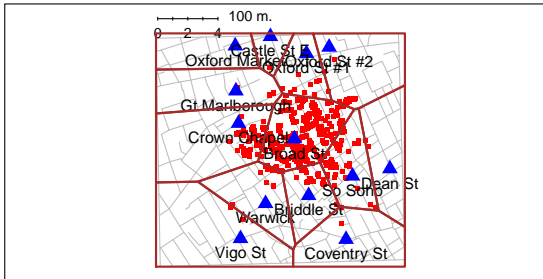


Figure 6.1: Deaths (red dots) and pump locations. Polygons that surround each pump enclose the locations for which that is the nearest pump.

Similar issues, arising from failures to ensure proper drainage systems, were repeated, from the 1840s and 1850s through until the end of the century, in New Zealand cities. See

Christine Dann, Sewage, water and waste - Stinking cities', Te Ara - the Encyclopedia of New Zealand, (8 June 2017)⁴

6.4 The Reinhoff and Rogoff saga (Smith, p.55)

Figure 6.2 plots data that underpinned the 2010 paper “*Growth in Time of Debt*” by the two Harvard economic historians Reinhart and Rogoff [RR]. The paper (Reinhart and Rogoff, 2010) has been widely quoted in support of economic austerity programs internationally. There was a huge stir, in the media and on the blogosphere, when graduate student Herndon found and published details of coding and other errors in the results that RR had presented.

```
##
## Attaching package: 'nlme'

## The following object is masked from 'package:HistData':
##
##      Wheat
```

⁴<https://teara.govt.nz/en/zoomify/24431/dunedin-renamed-stinkapool>

```
## This is mgcv 1.8-34. For overview type 'help("mgcv-package")'.
```

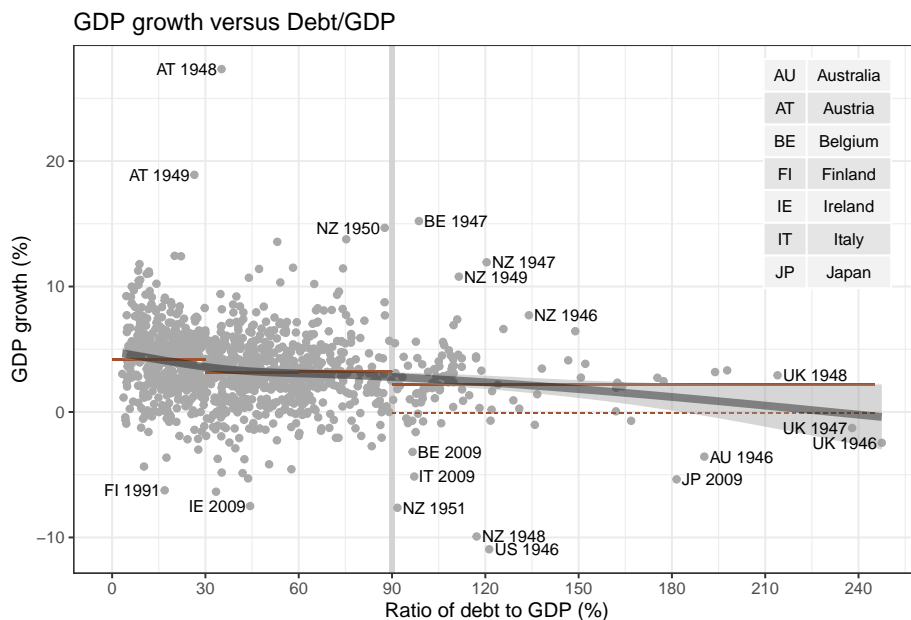


Figure 6.2: Red lines show means by Debt/GDP category. The means for $> 90\%$ Debt/GDP were from data for 10 only of the 20 countries. The red dotted line shows the (incorrect) mean given by RR. The smooth curve, fitted treating points as independent, shows no sudden change. It makes better sense to allow each country its own separate curve.

As well as coding errors, [Herndon](#)⁵ identified selective exclusion of available data, and unconventional weighting of summary statistics. There was a failure to acknowledge that the relationship studied has varied substantially by country and over time.

In response to [Herndon et al. \(2014\)](#) and other critics, Reinhart and Rogoff accepted that coding errors had led to the omission of several countries, but pushed back against other criticisms. Their revised analysis addresses only the

⁵http://www.peri.umass.edu/fileadmin/pdf/working_papers/working_papers_301-350/WP322.pdf

most egregious errors in their work. Among other issues, their insistence on treating each data point for each country as an independent piece of evidence makes no sense.

The smooth curve seems as good a summary as any, if points are, wrongly, to be treated as disconnected from country.

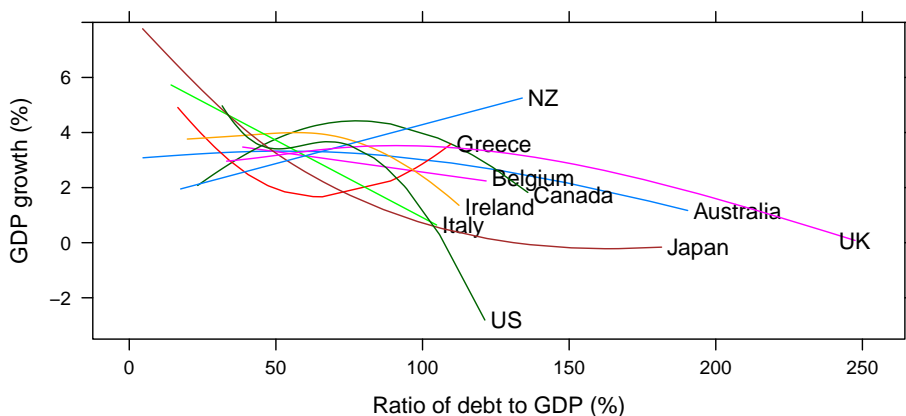


Figure 6.3: Smooths have been fitted for each of the 10 countries for which debt to GDP ratios were in some years greater than 90%. There is no consistent pattern, as there should be if RR's claim is to hold up.

Is there a pattern across countries?

```
av4 <- sapply(split(RR,RR$Country), function(x)
c(avT030 = mean(subset(x, debtgdp<=90)[ , 'dRGDP']),
av30T060 = mean(subset(x, debtgdp>30 & debtgdp<=60)[ , 'dRGDP']),
av60T090 = mean(subset(x, debtgdp>60 & debtgdp<=90)[ , 'dRGDP']),
av90plus = mean(subset(x, debtgdp>90)[ , 'dRGDP']),
nyears = nrow(subset(x, debtgdp>90))))
av4n <- av4[,!is.nan(av4[4,])]
dRGDPav <- sum(av4n[4,]*av4n[5,])/sum(av4[5,])
```

There were just 10 countries where debt to GDP ratios were greater than 90%

for one or more years. For 9 of those countries, the average of their GDP growth in the years at issue was on average positive relative to the previous year. The average percentage growth over all 10 countries, weighted according to number of years, was 2.168.

There are further serious issues of interpretation

- Does GDP drive debt/GDP ratio, or is it the other way round?
 - Or does a third factor drive both?
- Is the effect immediate, or on future economic performance

Smith (p.64) refers to work indicating that economic performance is more closely correlated with economic growth in the past than with future growth.

Parting comments

[Herndon et al. \(2014\)](#) comment

“... RR’s findings have served as an intellectual bulwark in support of austerity politics. The fact that RR’s findings are wrong should therefore lead us to reassess the austerity agenda itself in both Europe and the United States.”

The saga emphasizes the importance of working with reproducible code, rather than with spreadsheet calculations. The errors in RR’s calculations were from one perspective fortunate. Once highlighted, the errors drew critical attention to the paper, and to the serious flaws in the analysis.

Chapter 7

Weighting effects, & the Yule-Simpson Paradox

7.1 Deaths from Covid-19 — between country comparisons

United States data for the 13 months up to 31 January 2021 shows a stark difference in death rates, for those in younger as compared with those in older age groups, as shown in Figure 7.1.

```
cap_props <- paste0("United States reported deaths per 100 from Covid-19  
(and, for comparison, deaths per 100), for the time period from ",  
wkStartEnd[1], " to ", wkStartEnd[2], ". The total number  
reported was ", sum(USdeadpop1yr$covdeaths), ".")  
fromTo <- format(UScausesByCats[1,2:3], "%b %d %Y")
```

It is immediately obvious that more deaths are to expected, assuming a similar breakdown by age, in countries where there are more old people. Thus, the death rates per 1000 for Covid-19, for under age 65 as opposed to age ≥ 65 , were:

```
## <65 >=65  
## 0.28 6.30
```

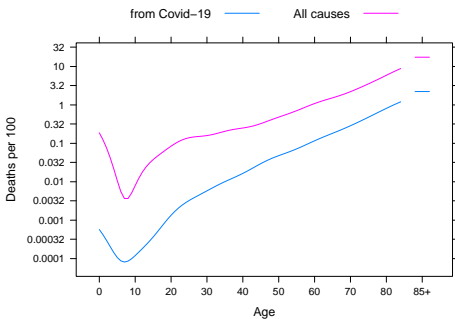


Figure 7.1: Proportions who died, from Covid-19, and in total

Now observe how these rates transfer across to countries with a different population structure.

| ## | US | Italy | Kenya |
|--------------------------------|-------|-------|-------|
| ## Percentage 65 or more | 16.3 | 23.3 | 2.5 |
| ## Expected deaths per 100,000 | 126.7 | 168.5 | 43.5 |
| ## Reported deaths per 100,000 | 126.7 | 146.0 | 3.3 |

Italy has a higher proportion of older people, leading to a higher expected death-rate. Kenya has a much lower proportion, and a lower expected overall death-rate. Other totals, e.g., for infection rates or hospital admissions, while still impacted by age structure, were not impacted to the same extent. A table on the US Center for Disease Control web page [Hospitalization and Death by Age] (<https://www.cdc.gov/coronavirus/2019-ncov/covid-data/investigations-discovery/hospitalization-death-by-age.html>) shows how, in the US, hospitalization and death varied with age. Whereas death rates for those 85 years of age or older were 810 times those of 5-17 year olds, for hospitalization the multiplier was 52. The rates for 0-4 year olds were around twice those for 5-17 year olds. Overall rates varied, in both cases, with ethnicity.

Between country comparisons can be hazardous. There are likely to be differences in the completeness of the data and in recording protocols. Case numbers are likely to be influenced by testing rates, are to be substantial undercounts, to an extent that varies from country to country.

7.2 The Yule-Simpson paradox — UCB admissions data

The Yule-Simpson paradox is a paradox of human intuition. It arises when weighting effects operate with respect to two or more factors. University of California Berkeley (UCB) 1973 admissions data for the six largest departments, summarized in Figure 7.2, provides an example. Admission rates varied by department, as did the relative numbers of males and females applying, in ways that led to the paradox. Male/female differences within departments were of lesser consequence.

Overall admission rates across the six largest departments were:

```
## Tabulate by Admit and Gender
byGender <- margin.table(UCBAdmissions, margin=1:2)
round(100*prop.table(byGender, margin=2)["Admitted", ], 1)
```

```
##      Male Female
##      44.5   30.4
```

Taken at face value, these numbers suggest discrimination against females,

```
par(mar=c(3.1,3.1,2.6,0.6), mgp=c(2,0.5,0))
applied <- margin.table(UCBAdmissions, margin=2:3)
pcAdmit <- 100*prop.table(UCBAdmissions, margin=2:3)["Admitted", , ]
byGender <- 100*prop.table(margin.table(UCBAdmissions,
                                         margin=1:2), margin=2)

dimnam <- dimnames(UCBAdmissions)
mfStats <- data.frame(Admit=c(pcAdmit[1,],pcAdmit[2,]),
                     Applicants=c(applied[1,], applied[2,]),
                     mf=factor(rep(dimnam[['Gender']],c(6,6)),
                                levels=dimnam[['Gender']]),
                     Department=rep(dimnam[["Dept"]],2))

xlim <- c(0, max(mfStats$Admit)*1.025)
ylim <- c(0, max(mfStats$Applicants)*1.075)
plot(Applicants~Admit, data=mfStats, xlim=xlim, ylim=ylim,
     fg="gray", cex.lab=1.25,
     col=rep(c("blue","red"),rep(6,2)),type="h",lwd=2,
     xlab="UCB Admission rates (%) , 1973",
     ylab="Number of applicants")
```

```

pcA <- rbind(pcAdmit[1,], apply(pcAdmit,2, mean)+2, pcAdmit[2,], rep(NA,6))
pcA[2,3] <- pcA[2,3]+1
appA <- rbind(applied[1,], apply(applied,2, mean)+80,
              applied[2,], rep(NA,6))
deptNam <- dimnam[[3]]
for(j in 1:ncol(appA))
  lines(pcA[,j], appA[,j], type="c", col="gray")
text(pcA[2,],appA[2,],deptNam)
##
par(xpd=TRUE)
text(byGender[1,1:2], rep(par()$usr[4],2)+0.5*strheight("^"), labels=c("^", "^"),
     col=c("blue", "red"), cex=1.25, srt=180)
text(byGender[1,], par()$usr[4]+1.4*strheight("A"),
     labels=paste(round(byGender[1,],1)), cex=0.85)
text(byGender[1,1:2]+c(-3.5,3.5), rep(par()$usr[4],2)+2.65*strheight("A"),
     labels=c("All males", "All females"), pos=c(4,2), cex=1.25)
par(xpd=FALSE)
abline(h=200*(0:4), col="lightgray", lty="dotted")
abline(v=20*(0:4), col="lightgray", lty="dotted")
legend("topleft", col=c('blue', 'red'), lty=c(1,1), lwd=0.75, cex=0.9,
      y.intersp=0.65, legend=c("Males", "Females"), bty="n")

```

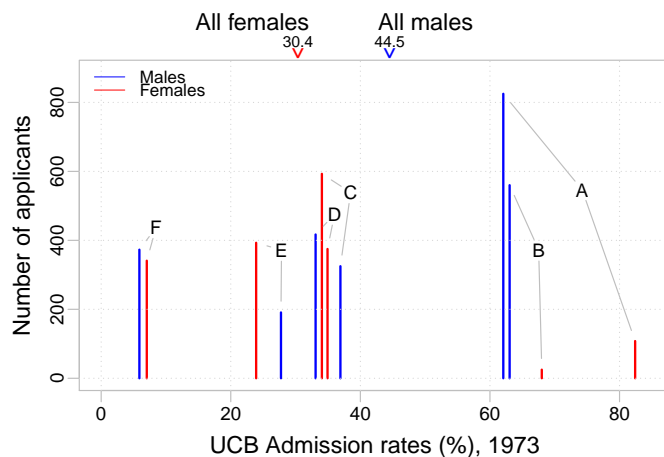


Figure 7.2: UCB admission data for 1973, accounting for male/female differences, by department. Department labels range from A to E.

Admission rates for males and females, by department, tell a different story:

```
admit <- round(100*prop.table(UCBAdmissions,
                              margin=2:3)["Admitted", , ], 1)
tab <- xtable::xtable(admit, digits=1)
print(tab, type='latex', comment=FALSE, floating=FALSE, scalebox=0.9)
```

| | A | B | C | D | E | F |
|--------|------|------|------|------|------|-----|
| Male | 62.1 | 63.0 | 36.9 | 33.1 | 27.7 | 5.9 |
| Female | 82.4 | 68.0 | 34.1 | 34.9 | 23.9 | 7.0 |

The biggest differences in admission rates were in departments A (82.4%-62.1%=20.3%) and B (68%-63%=5%), in both cases favouring females. In the other four departments, differences were 3.8% or less.

In order to understand how this happens, it is necessary to look at how the numbers that applied break down by gender and by department:

```
library(xtable)
tots <- margin.table(UCBAdmissions, margin=2:3)
tab <- xtable::xtable(tots, digits=0)
print(tab, type='latex', comment=FALSE, floating=FALSE, scalebox=0.9)
```

| | A | B | C | D | E | F |
|--------|-----|-----|-----|-----|-----|-----|
| Male | 825 | 560 | 325 | 417 | 191 | 373 |
| Female | 108 | 25 | 593 | 375 | 393 | 341 |

Comparing departments A and B with other departments, one finds:

```
##
##               AB               CDEF
## Male   1385 (62.1,63%)  981 (5.9 to 36.9%)
## Female  133 (82.4,68%) 1109 (7.0 to 34.9%)
```

- Overall admission rates for males are weighted (1385:981) towards male rates of 62.1% or 63% in departments A and B.
- Those for females are strongly weighted (1109:133) towards admission rates that range from 5.9% to 34.1% in departments C to F.

UCB Admissions Data – Another perspective

```
par(mar=c(3.1,3.1,2.1,0.6), mgp=c(2,0.5,0))
diffAdmit <- t(pcAdmit)
diffAdmit <- data.frame(dept=rownames(diffAdmit), stringsAsFactors=FALSE,
                        Male=diffAdmit[,1], Female=diffAdmit[,2],
                        diff=diffAdmit[,2]-diffAdmit[,1])
diffAdmit$mfratio <- signif(apply(applied,2, function(x)x[1]/x[2]),2)
ord <- order(diffAdmit$diff)
plot(diff ~ I(1:6), data=diffAdmit[ord,], pch=dept, cex=1.25,
     ylim=c(-5.25,21.5), xlab="Order of difference", fg="gray",
     ylab="Female - Male difference (%)")
with(diffAdmit[ord,], text(diff~I(1:6), labels=paste(mfratio),
              pos=c(rep(4,5),2)), cex=0.8)
leg <- c("Male:Female ratios of numbers",
        "who applied are shown alongside",
        "each letter")
legend("topleft", legend=leg, box.col="gray")
```

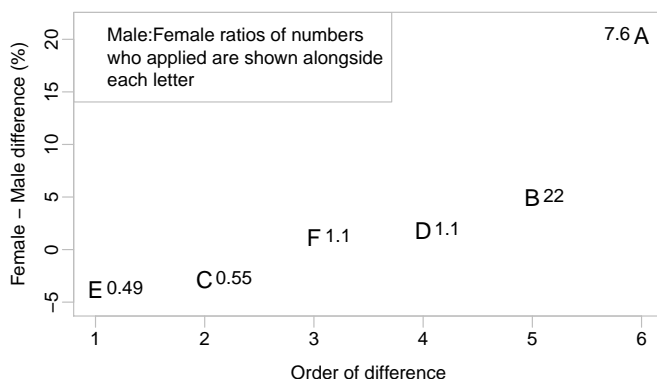


Figure 7.3: UCB admission data for 1973 — another perspective.

Effective graphs can help greatly in fairly representing the evidence. The website <https://www.youtube.com/watch?v=ZDinnCwP3dg> has an animated video that provides a short overview of the paradox.

The same sorts of paradoxical effects can be found in regression. The Yule-Simpson paradox may be regarded as a special case of Lord's paradox, described in Lord (1967). Any attempt to attach meaning to regression coefficients can be highly misleading, unless one can be sure that effects of all variates and covariates are properly accounted for. It is rarely easy, with observational data, to be sure that this has been done effectively.

7.3 Third variables change the story — further examples

Here, “variable” is used as generic name for either continuous variables, or categorical variables, otherwise called factors.

Does Baclofen help in reducing pain?

```

cap9 <- "Data are pain reduction scores. Subgroup numbers, shown
        below each point, weight the overall average ('ALL') for
        baclofen towards the high female average, and for no baclofen
        slightly towards the low male average."

library(lattice)
parset <- simpleTheme(cex=1.35, pch=16,
                      col=c("darkblue", "turquoise"))
gabalong <- data.frame(values=unlist(DAAG::gaba["30",])[-1],
                      sex=rep(c("male", "female", "ALL"), rep(2,3)),
                      trt=rep(c("Baclofen", "No baclofen"), 3))
gph <- stripplot(sex~values, groups=trt, data=gabalong,
                 par.settings=parset,
                 xlab=list("Average reduction: 30 min vs 0 min",
                           cex=1.0),
                 scales=list(cex=1.0),
                 panel=function(x,y,...){
                   panel.stripplot(x,y,...)
                   ltext(x,y,paste(c(3,9,15,7,22,12)), pos=1,
                          cex=0.8)
                 }, auto.key=list(columns=2, points=TRUE, cex=1.0))
gph + latticeExtra::layer(panel.abline(h=1.35, col="gray"))

```

Researchers were comparing two analgesic treatments, without and with baclofen. When the paper was first submitted for publication, an alert reviewer spotted that some of the treatment groups contained more women than men, and asked whether this might account for the results.¹

For a fair overall comparison:

- Calculate means for each subgroup separately.
- Overall treatment effect is average of subgroup differences.

The effect of baclofen (reduction in pain score from time 0) is then:

- Females: $3.479 - 4.151 = -0.672$ (-ve, therefore an increase)
- Males: $1.311 - 1.647 = -0.336$
- Average, male & female = $-0.5 \times (0.672 + 0.336) = -0.504$

¹Cohen, P. 1996. Pain discriminates between the sexes. *New Scientist*, 2 November, p. 16.

7.3. THIRD VARIABLES CHANGE THE STORY — FURTHER EXAMPLES57

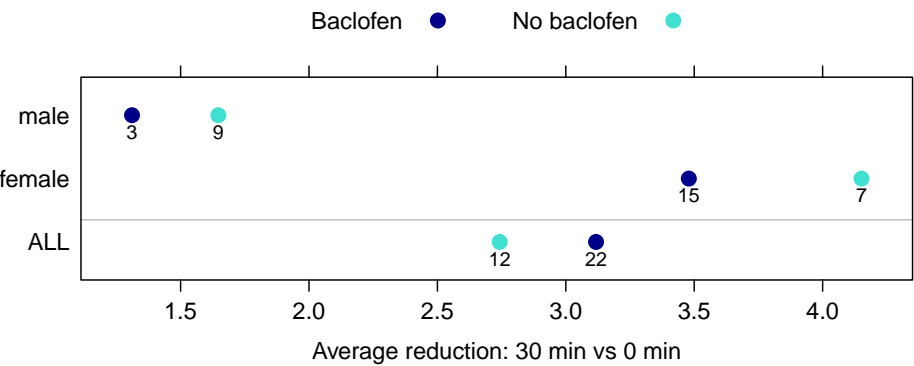


Figure 7.4: Data are pain reduction scores. Subgroup numbers, shown below each point, weight the overall average (“ALL”) for baclofen towards the high female average, and for no baclofen slightly towards the low male average.

Web page revenue per click (Smith, p.111)

In 2010 a US Internet company collected data on the effectiveness of two different web page layouts:

- 1-click: Adverts appear on website’s first page
- 2-click: Click on keyword sends user to page with advert
 - Check interest, then show advert

| 1-click | | | 2-click | | |
|---------|-------|---------|---------|-------|---------|
| Revenue | Users | RPM | Revenue | Users | RPM |
| \$2.9 | 250 | \$11.60 | \$1.7 | 140 | \$12.14 |

- Revenue and users are millions
- RPM is revenue per million users

The data combines US and International users. Is that an issue?

US vs International users

| Location | 1-click | | | 2-click | | |
|----------|---------|-----------------------|---------|---------|---------------------|---------|
| | Revenue | Users | RPM | Revenue | Users | RPM |
| US | \$1.8 | 70 | \$27.71 | \$1.2 | 50 | \$24.00 |
| Int. | \$1.1 | 180 | \$6.11 | \$0.50 | 90 | \$5.56 |
| | | $\frac{180}{70}=2.57$ | | | $\frac{90}{50}=1.8$ | |

- Revenue and users are millions
- RPM is revenue per million users

This was followed up with a randomized experiment (an A/B test.)

7.4 Cricket Bowling Averages

Runs (R), wickets (W) and runs per wicket ($\{RPW\}$)

| | 1st innings | | | 2nd innings | | | Overall | | |
|----------|-------------|---|-------|-------------|---|-------|---------|----|-------|
| | R | W | RPW | R | W | RPW | R | W | RPW |
| Bowler A | 40 | 4 | 10.0 | 240 | 6 | 40.0 | 280 | 10 | 28.0 |
| Bowler B | 70 | 5 | 14.0 | 50 | 1 | 50.0 | 120 | 6 | 20.0 |

Fair comparison: Compare runs per wicket ($\{RPW\}$)

| | 1st innings | | 2nd innings | | Overall | |
|----------|-------------|-----|-------------|-----|------------------------|------|
| | RPW | W | RPW | W | RPW | W |
| Bowler A | 10.0 | (4) | 40.0 | (6) | $\frac{10+40}{2} = 25$ | (10) |
| Bowler B | 14.0 | (5) | 50.0 | (1) | $\frac{50+14}{2} = 32$ | (6) |

7.5 Epistatic effects on genetic studies

In population genetics, Simpson's paradox type effects are known as epistasis. Most human societies are genetically heterogeneous. In San Francisco, any gene that is different between the European and Chinese populations will be found to be associated with the use of chopsticks! If a disease differs in frequency between the European and Chinese populations, then a naive analysis will find an association between that disease and any gene that differs in frequency between the European and Chinese populations.

Such effects are major issues for gene/disease population association studies. It is now common to collect genetic fingerprinting data that should identify major heterogeneity. Providing such differences are accounted for, large effects that show up in large studies are likely to be real. Small effects may well be epistatic.

Chapter 8

Regression and Correlation

Yule-Simpson type effects, discussed in Section 7 [“Explaining the Yule-Simpson Paradox”], are important in a regression context also. Nonsense correlations that arise where the third variable is time provide simple examples.

8.1 What direction does the correlation go?

Variable A may cause variable B. Or variable B may cause variable A. Or both A and B may be caused (or driven) by a third variable C. Cases where the third variable is time, as in Figure 8.1, are a fruitful source of examples of spurious correlations.¹

```
cap11 <- "Sociology PhDs awarded (from US National Science
Foundation data) vs Deaths from Anticoagulants.
"
```

```
suppressPackageStartupMessages(library(latticeExtra))
phdDeaths <- data.frame(Year = 1999:2009,
Doctorates = c(572,617,566,547,597,580,536,579,576,601,664),
Deaths = c(17,39,39,27,44,46,29,42,47,52,78))
```

¹Figure 8.1 is one of many such examples that are available from <http://www.tylervigen.com/spurious-correlations>.

```
obj1 <- xyplot(Doctorates ~ Year, data=phdDeaths, type="b", xlab="")
obj2 <- xyplot(Deaths ~ Year, data=phdDeaths, type="b")
update(doubleYScale(obj1, obj2, add.ylab2=TRUE),
       par.settings=simpleTheme(col=c("purple", "black"),
                                pch=16, cex=2))
```

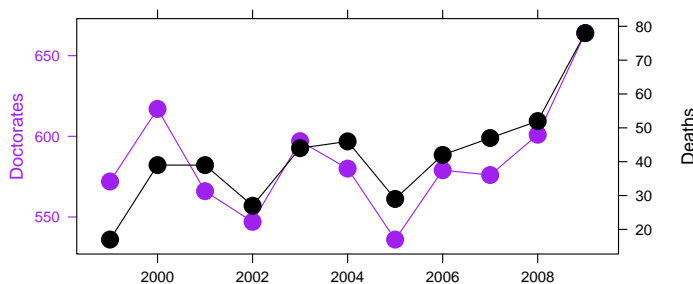


Figure 8.1: Sociology PhDs awarded (from US National Science Foundation data) vs Deaths from Anticoagulants.

Examples of this type help highlight how correlation can and cannot reasonably be used.

The following examples, where third variables are likely to be involved, are from (Nisbett, 2016, pages 133-134).

1. Children of parents who try to control eating are more likely to be overweight.
2. Countries with higher IQs have higher average wealth measures.
3. People who smoke marijuana are more likely to later use cocaine.
4. Ice cream consumption & polio were closely correlated in the 1950s.

8.2 Regression to the mean

Tall fathers are likely to have tall sons, but shorter than themselves. Tall sons are likely to have tall fathers, but shorter than themselves. The data shown in Figure 8.2 are from Karl Pearson. See `?HistData::PearsonLee`.

cap12 <- "Tall fathers are likely to have tall sons, but shorter than themselves.
Tall sons are likely to have tall fathers, but shorter than themselves."

```
par(mar=c(3.6,3.6,2.1, 1.1), mgp=c(2,0.5,0))
ptCol <- adjustcolor('black', alpha=0.4)
load('Data/Pearson.RData')
ptCol <- adjustcolor('black', alpha=0.4)
plot(Pearson[,1:2], col=ptCol, xlab="Father's height (in)", ylab="Son's height (in)")
avSon <- mean(Pearson$Son)
avFather <- mean(Pearson$Father)
points(avFather, avSon, cex=1.5, col=2, pch=16)
r <- cor(Pearson)[1,2]
text(60.5,77, paste("r =", round(r,2)))
Pearson.lm <- lm(Son~Father, data=Pearson)
b <- coef(Pearson.lm)
pc95 <- quantile(Pearson$Father, 0.95)
son95 <- quantile(Pearson$Son, 0.95)
pred95 <- b[1]+b[2]*pc95
abline(Pearson.lm, col=2)
abline(v=pc95, col=2)
abline(h=pred95, col=2)
abline(h=son95, col=2, lty=2)
pcSon <- round(100*with(Pearson, sum(Son<pred95)/nrow(Pearson)))
mtext(side=3, line=0.05, "95th percentile", col=2, at=pc95)
mtext(side=4, line=0.85, paste0(pcSon,"th\npercentile"), col=2, at=pred95, srt=90)
mtext(side=3, line=1.0, "Son's height, given father's height (red lines)")
plot(Pearson[,1:2], col=ptCol)
text(60.5,77, paste("r =", round(r,2)))
Pearson.lm2 <- lm(Father~Son, data=Pearson)
c2 <- coef(Pearson.lm2)
b2 <- c(-c2[1]/c2[2], 1/c2[2])
son95 <- quantile(Pearson$Son, 0.95)
predF95 <- c2[1]+c2[2]*son95
abline(Pearson.lm, col=adjustcolor("red",alpha=0.4))
abline(b2, col="purple", lty=2, lwd=2)
abline(h=son95, col="purple", lty=2, lwd=2)
abline(v=predF95, col="purple", lty=2, lwd=2)
abline(v=pc95, col=adjustcolor('red',alpha=0.4))
```

```
abline(h=pred95, col=adjustcolor('red',alpha=0.4))
pcF <- round(100*with(Pearson, sum(Father<=predF95)/nrow(Pearson)))
mtext(side=3, line=1.0, "Father's height, given son's height (purple lines)")
mtext(side=3, line=-0.25, paste0(pcF,"th percentile"), col="purple", at=predF95)
mtext(side=4, line=0.05, paste0("95th percentile"), col="purple", at=son95, sr
```

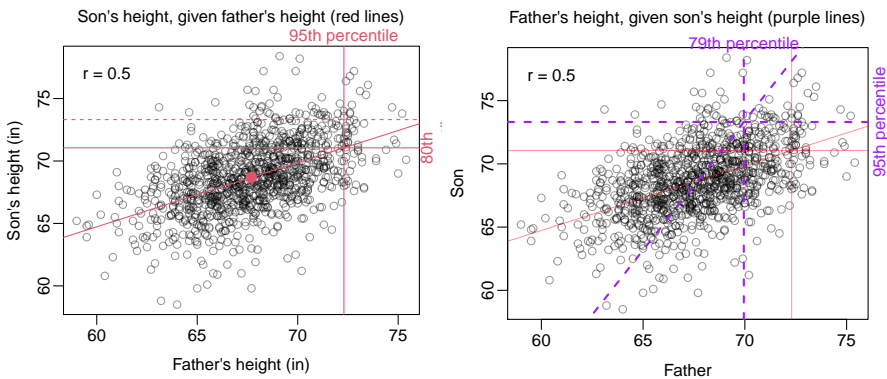


Figure 8.2: Tall fathers are likely to have tall sons, but shorter than themselves. Tall sons are likely to have tall fathers, but shorter than themselves.

Kahneman argues, perhaps too simplistically:

- Height is mainly due to genetic factors
- Sons share half their genes with their fathers
- Hence, correlation between sons' & fathers' heights $\simeq 0.5$

Galton's 1886 data, which predates Pearson's data, shows a 0.46 correlation between child height and the average of the parent height.

Regression to the mean in verse:

<https://www.youtube.com/watch?v=sxMlckUWaw>

Kahneman's comments on regression to the mean

"Extreme predictions and a willingness to predict rare events from weak evidence are both manifestations of System 1. ..."

"Regression to the mean is also a problem for System 2. The very idea ... is alien and difficult to communicate and comprehend. This is a case where System 2 requires special training."

“We intuitively want to match predictions to the evidence.”

“We will not learn to understand regression from experience.”

8.3 Regression to the mean in a variety of contexts

Decathlon scores — between event correlations

```
cap13 <- "Between event correlations for top performances in the decathlon
between 1985 and 2006."

Decath2006 <- subset(GDadata::Decathlon,yearEvent==2006)
pointsNam <- c("100m", "Longjmp", "Shotput", "Highjmp", "400m", "110mHurd",
"Discus", "PVault", "Javelin", "1500m")
## ord <- dput(corrplot::corrMatOrder(cor(Decath2006[,15:24])))
ord <- c(10L, 9L, 8L, 3L, 7L, 4L, 6L, 2L, 1L, 5L)
## Check difference of rank from Pearson correlations
corP <- cor(Decath2006[,ord+14])
corS <- cor(Decath2006[,ord+14],method="spearman")
## round(corP-corS,2)
## dput(corrplot::corrMatOrder(cor(Decath2006[,15:24])))
ord <- c(10L, 9L, 8L, 3L, 7L, 4L, 6L, 2L, 1L, 5L)

##
## Attaching package: 'ggplot2'

## The following object is masked from 'package:latticeExtra':
##
##     layer

## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg     ggplot2

my_fn <- function(data, mapping, pts=list(), smt=list(), ...){
  ggplot(data = data, mapping = mapping, ...) +
    do.call(geom_point, pts) +
```

```

do.call(geom_smooth, smt)
}
my_custom_cor <- function(data, mapping, color = I("grey50"), sizeRange = c(1,

  # get the x and y data to use the other code
  x <- eval_data_col(data, mapping$x)
  y <- eval_data_col(data, mapping$y)
  ct <- cor.test(x,y)

  r <- unname(ct$estimate)
  rt <- format(r, digits=2)[1]

  # since we can't print it to get the strsize, just use the max size range
  cex <- max(sizeRange)

  # helper function to calculate a useable size
  percent_of_range <- function(percent, range) {
    percent * diff(range) + min(range, na.rm = TRUE)
  }

  # plot the cor value
  ggally_text(
    label = as.character(rt),
    mapping = aes(),
    xP = 0.5, yP = 0.5,
    size = I(percent_of_range(cex * abs(r), sizeRange)),
    color = color,
    ...
  ) +
  # remove all the background stuff and wrap it with a dashed line
  theme_classic() +
  theme(
    panel.background = element_rect(
      color = color,
      linetype = "longdash"
    ),
    axis.line = element_blank(),
    axis.ticks = element_blank(),

```

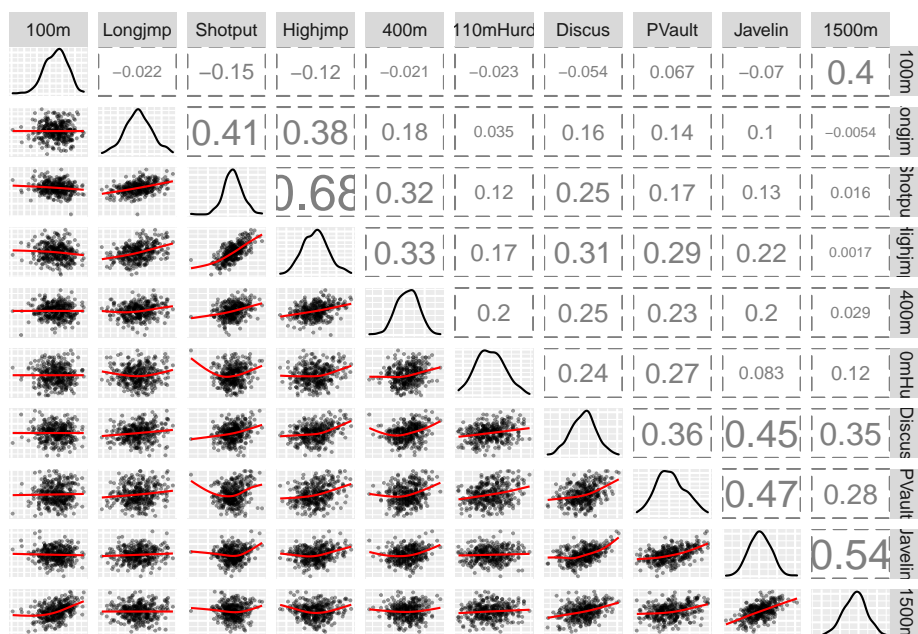


```

    axis.text.y = element_blank(),
    axis.text.x = element_blank()
  )
}

ggpairs(Decath2006[, (15:24)[ord]],
        lower = list(continuous =
                      wrap(my_fn,
                           pts=list(size=0.25, colour=adjustcolor("black", alpha=0.5),
                                   smt=list(method="gam", se=F, size=0.5, colour="red"))),
                      axisLabels="none", columnLabels=pointsNam,
                      upper=list(continuous=wrap(my_custom_cor, sizeRange=c(2.5,4.5))))

```

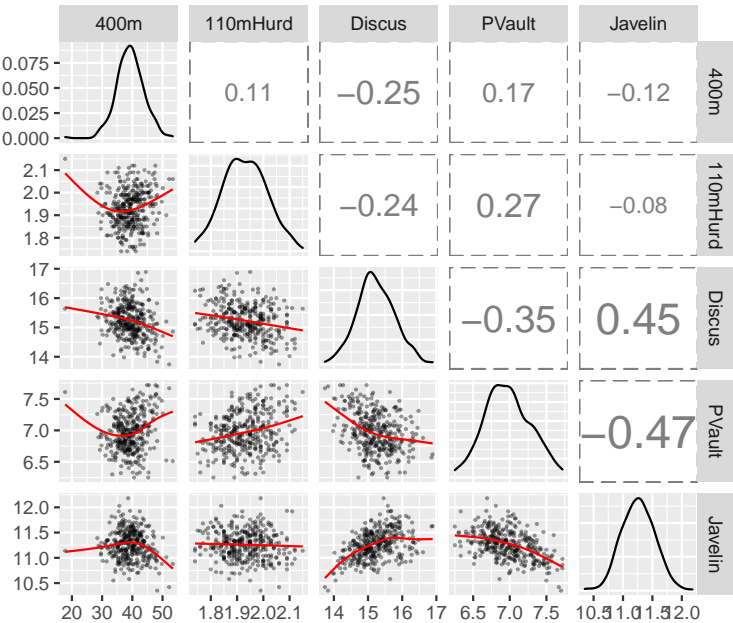


The dataset `GDadata::Decathlon` has best annual performances of 6800 points and over for a twenty-one year period after new rules were introduced in 1985.²

²Data are from the Estonian website <http://www.decathlon2000.com/>

Plot measures, not points (selection only)

```
choosecols <- ord[5:9]
ggpairs(Decath2006[, (4:13)[choosecols]],
        lower = list(continuous =
                      wrap(my_fn,
                           pts=list(size=0.25, colour=adjustcolor("black", al
                           smt=list(method="gam", se=F, size=0.5, colour="red
        columnLabels=pointsNam[5:9],
        upper=list(continuous=wrap(my_custom_cor, sizeRange=c(3,5))))
```



Total profit to total income ratio, by industry class

Data are from the Statistics NZ Business Performance Benchmark³

³<https://shinyapps.stats.govt.nz/bpb/>

```

cap14 <- paste("Total profit to total income ratio, by industry class ---",
"correlation between 2015 value and 2013 value.")

load('data/FinAllInd.RData')
expend <- subset(FinAllInd, subset=Variable=="Expenditure")
profit <- subset(FinAllInd, subset=Variable=="Profit")
income <- subset(FinAllInd, subset=Variable=="Income")
assets <- subset(FinAllInd, subset=Variable=="Assets")
employ <- subset(FinAllInd, subset=Variable=="Employee Count")
# expend[,6] <- as.numeric(expend[[6]])
# expend[,7] <- as.numeric(expend[[7]])
# expend[,8] <- as.numeric(expend[[8]])
plot(I(profit[[9]]/income[[9]])~I(profit[[11]]/income[[11]]),
     pch=16, cex=1.25, xlab="2013 ratio", ylab="2015 ratio",
     cex.lab=1.55, col=adjustcolor('purple',alpha=0.5), fg="gray",
     xlim=c(-0.15,0.475), ylim=c(-0.125, 0.41))
mtext(side=3, line=0.5, "Total profit to total income ratio")
par(xpd=TRUE)
rt <- c(31,37,167, 330, 385, 441, 458)
left <- c(10, 15, 35, 57, 76,122,157,175, 164, 196, 253, 274, 292,341,362, 384,430)
halfblack <- adjustcolor('black', alpha=0.35)
df <- data.frame(ratio15=I(profit[[9]]/income[[9]]), ratio13=I(profit[[11]]/income[
indclass=profit[["Industry_class"]])
with(df[rt,], text(ratio15~ratio13,col=halfblack,labels=indclass,pos=2))
with(df[left,], text(ratio15~ratio13,col=halfblack,labels=indclass,pos=4))

```

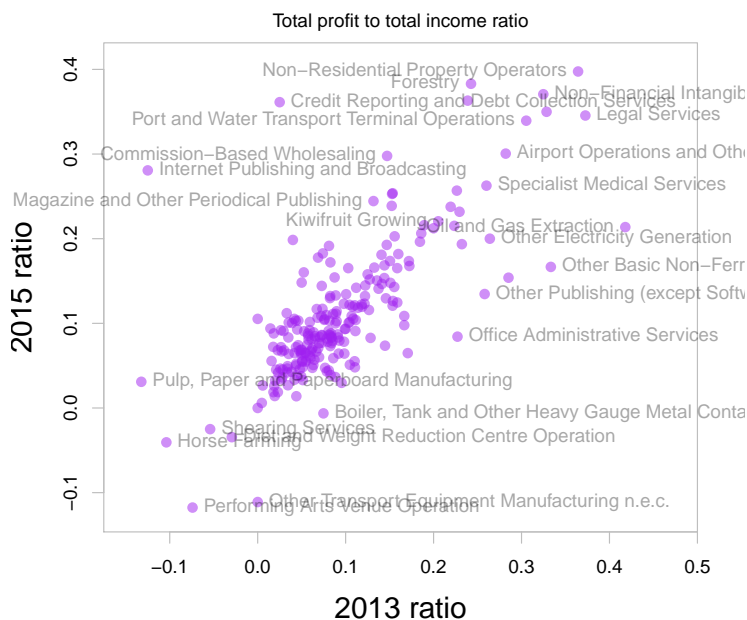


Figure 8.3: Total profit to total income ratio, by industry class — correlation between 2015 value and 2013 value.

Not shown are 4 points where the 2015 ratio was less than -0.3.
(Olive growing, Petroleum exploration, Stevedoring Services.)

NBA player total points — correlations decline over time

```
cap12.5 <- "As time progresses, correlation decreases, and regression to the m

load("data/NBAplayer.RData")
m1516 <- with(NBAplayer, match(z15$Name,z16$Name,nomatch=0))
df <- with(NBAplayer, data.frame(p1516=z15$TotalPoints[m1516>0], p1617=z16$Tot
xyplot(p1617~p1516, data=df,
      xlab="Total points, 2015-2016",
      ylab="Total points, 2016-2017", type=c("p","r"),
      main=list("A: 2016-2017 vs 2015-2016", font=1, y=0))
```

```

m1116 <- with(NBAplayer, match(z11$Name,z16$Name,nomatch=0))
df <- with(NBAplayer, data.frame(p1112=z11$TotalPoints[m1116>0], p1617=z16$TotalPoints[m1617>0]))
xyplot(p1617~p1112, data=df,
       xlab="Total points, 2010-2011",
       ylab="Total points, 2016-2017", type=c("p","r"),
       main=list("B: 2016-2017 vs 2010-2011",font=1, y=0))

```

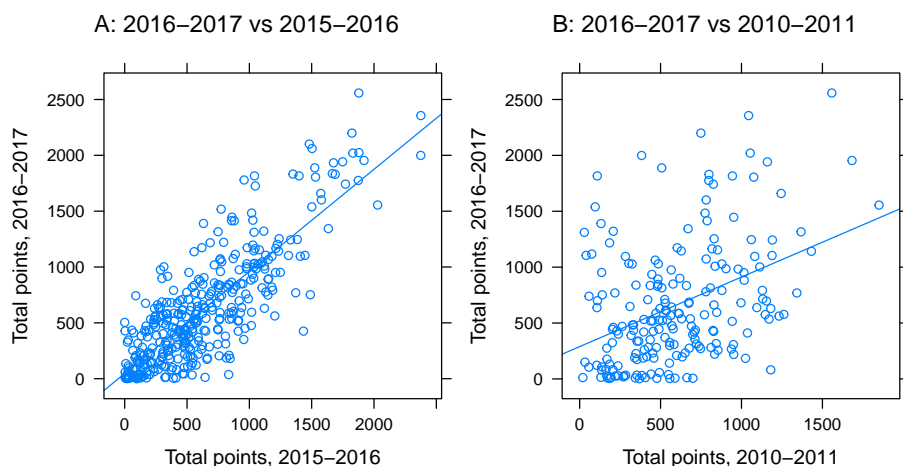


Figure 8.4: As time progresses, correlation decreases, and regression to the mean increases.

The Sports Illustrated cover “jinx”

The 21 January 2002 issue of Sports Illustrated featured an [article on the so-called The Sports Illustrated cover jinx] (<https://vault.si.com/vault/2002/01/21/that-old-black-magic-millions-of-superstitious-readersand-many-athletesbelieve-that-an-appearance-on-sports-illustrateds-cover-is-the-kiss-of-death-but-is-there-really-such-a-thing-as-the-si-jinx>)⁴ After examining covers that had appeared from 1954 to 2001, they found that

⁴<https://vault.si.com/vault/2002/01/21/that-old-black-magic-millions-of-superstitious-readersand-many-athletesbelieve-that-an-appearance-on-sports-illustrateds-cover-is-the-kiss-of-death-but-is-there-really-such-a-thing-as-the-si-jinx>

Of the 2,456 covers SI had run, 913 featured a person who, or team that, suffered some verifiable misfortune that conformed to our definition—a Jinx rate of 37.2%.

[Wikipedia](#) has very detailed comments on the phenomenon. Especially in contact sports, athletes who do exceptionally well will have been pushing their bodies to the limit, with a high risk of injury.

Most athletes that seemed to suffer the jinx most typically suffered because of an injury to their body, or some other bad luck following their appearance.

Athletes appear on the cover of a magazine such as “Sports Illustrated” when they are performing unusually well, both relative to their fellow athletes, and relative to their own earlier performances. They are likely to be approaching their peak, and/or to have experienced an unusual run of luck. Where chance effects are the main driver, their success is likely to be temporary, quickly dropping back to close to or below their longer term average. Some of those who feature will be at their peak, and will on that account soon drop back.

The Wikipedia article notes a number of athletes who went on to do exceptionally well following appearance on the cover. Standouts were the basketballer Michael Jordan (50 covers) and Muhammad Ali (40 covers). These were outstanding athletes that consistently outperformed others.

8.4 Moderating subjective assessments

- Estimate average
- Make assessment, based on what evidence seems to suggest
 - Assessment is based on current measure
- Estimate correlation between current and predicted measure.
- The correlation determines the fraction of the distance to move from the average to the assessment.

e.g., Correlation of shotput with long jump is close to 0.4

- 14.9 meter put is at the 93nd centile (7% will do better)
- Long jump mean = 6.97; 92% mark 7.47 (difference=0.5)
- Estimate long jump result as $6.97 + 0.4 \times 0.5 = 7.17$

Forecasting sales

You are the sales forecaster for a department store chain.

All stores are similar in size and merchandise offered, but random factors affect sales in any year. Overall sales are expected to increase by 10% from 2020 to 2021. Sales in 2020, with the expected total and mean for 2021 are, in millions of dollars:

| Store | 2020 | 2021 |
|-------|------|------|
| 1 | 10 | — |
| 2 | 23 | — |
| 3 | 18 | — |
| 4 | 29 | — |
| TOTAL | 80 | 88 |
| MEAN | 20 | 22 |

Assuming that year to year correlation is around 0.4, what is a reasonable estimate for sales in each of the stores in 2021. The mean sales amount in 2021 is predicted to be 22,000,000 dollars? With a correlation of 0.4, the predicted sales for the individual stores are obtained thus:

| Store | 2020 | Subtract 20 | Xply by 0.4, add to 22 | Predicted sales |
|-------|------|-------------|------------------------|-----------------|
| 1 | 10 | -10 | 22-4 | 18 |
| 2 | 23 | +3 | 22+1.2 | 23.2 |
| 3 | 18 | -2 | 22-0.8 | 21.2 |
| 4 | 29 | +9 | 29+3.6 | 32.6 |
| MEAN | 20 | 0 | 22 | 22 |

Choosing from job applicants

Correlation between presentation & performance is likely to be lower for the less well-known. In both cases performance is likely, relative to presentation, to move in closer to the mean. For less well-known candidates, the shift towards the mean is likely to be greater.

Secrist's "The Triumph of Mediocrity in Business"

Horace Secrist's 1933 book was based on annual data for 1920 to 1930:

- 73 different industries; examine ratios
 - Profits:sales; Profits:assets; Expenses:sales; Expenses:assets
- For each industry in 1920: split firms into 4 quartiles: top 25%, ...
 - Took average for each statistic, for each quartile, for each year.
 - Surprise, surprise, the best went, on average, down ...

"Complete freedom to enter trade and the continuance of competition mean the perpetuation of mediocrity. ... neither superiority or inferiority will tend to persist. Rather, mediocrity tends to become the rule."

Note again — regression to the mean goes in both directions!

```
cap15 <- paste("Simulations based on correlation for Secrist's data ---
               showing how regression to the mean goes in both directions.")

par(mfrow=c(1,2), mgp=c(2,0.5,0), mar=c(3.1,3.1,1.6,1.1))
set.seed(29)
x1 <- rnorm(1000, mean=34, sd=5)
x2 <- 0.7*x1 + rnorm(1000, mean=.3*34, sd=5*sqrt(1-.49))
x3 <- 0.7*x2 + rnorm(1000, mean=.3*34, sd=5*sqrt(1-.49))
## print(round(c(sd(x1),sd(x2),sd(x3)),2))
## print(round(c(range(x1),range(x2),range(x3)),2))
xx <- data.frame(x1=x1,x2=x2,x3=x3)
xx <- xx[order(xx[,1]),]
qxx <- sapply(xx, function(x)c(mean(x[1:250]), mean(x[251:500]),
  mean(x[501:750]),mean(x[751:1000])))
plot(c(1920,1925,1930),qxx[4,], xlab="", ylab="Profit",
      ylim=c(26.6,40.7), fg="gray", type="n")
for(i in 1:4){
  lines(c(1920,1925,1930),qxx[i,], ylab="Profit", fg="gray", type="b", pch=16)
}
text(rep(1920,4), qxx[,1], c("Lowest 25% in 1920","2nd Lowest 25% in 1920",
  "2nd highest 25% in 1920", "Highest 25% in 1920"))
```



```

pos=4, col=adjustcolor("blue",alpha=0.5))
xx2 <- xx[order(xx[,3]),]
qxx2 <- sapply(xx2, function(x)c(mean(x[1:250]), mean(x[251:500]),
  mean(x[501:750]),mean(x[751:1000])))
plot(c(1920,1925,1930),qxx[4,], xlab="", ylab="Profit",
  ylim=c(26.6,40.7), fg="gray", type="n")
for(i in 1:4){
lines(c(1920,1925,1930),qxx2[i,], ylab="Profit", fg="gray", type="b", pch=16)
}
text(rep(1930,4), qxx2[,3], c("Lowest 25% in 1930","2nd Lowest 25% in 1930",
  "2nd highest 25% in 1930", "Highest 25% in 1930"),
  pos=2, col=adjustcolor("blue",alpha=0.5))

```

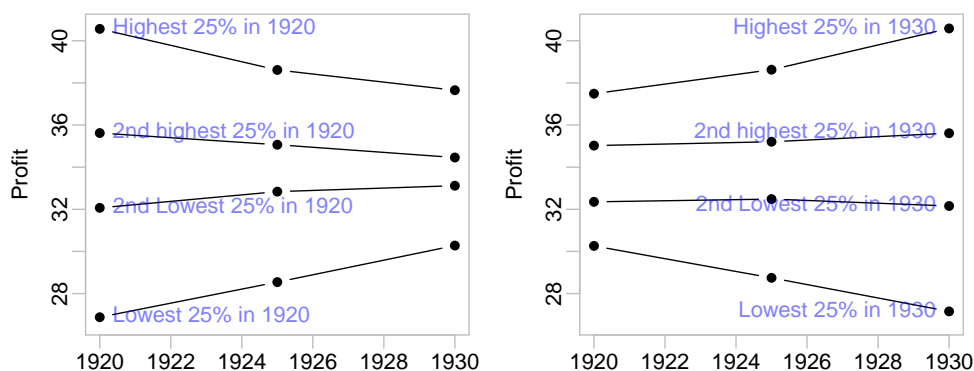


Figure 8.5: Simulations based on correlation for Secrist’s data — showing how regression to the mean goes in both directions.

“Do old fallacies ever die?”

Smith gives references to work by prominent economists in the past half-century that had quoted Secrist approvingly or repeated his error.

- 1970: “The book [by Secrist] contains an elaborate statistical demonstration that, over a period of time, initially high-performing ...”
- 1980s investment textbook: “Ultimately, economic forces will force the

convergence of the profitability and growth rates of different firms.” This was backed up with a 1980/1966 Secrist type comparison.

- 2000: (Journal article) “... profitability is mean-reverting within as well as across industries. Other firms eventually mimic products and technologies that produce above normal profitability ...”
- Friedman (2003) “Do old fallacies ever die?” cites other examples.

Chapter 9

Covariate adjustments in observational studies

At least in principle, it is relatively straightforward to use regression type methods to make predictions for a set of new data that have been sampled in the same way. What is hard for observational data, much harder than is commonly acknowledged, is to give the model coefficients a causal interpretation. For this, it is necessary to have a clear understanding of the processes involved.

- There will be several, perhaps a very large number, of explanatory variables, and an outcome variable.
- The aim is to find a model that will make predictions for new data.
- Note the predictive/descriptive distinction.
 - How do engineers predict building risk?
 - Note the “in sample/out of sample” distinction.
 - But is the “new” a random sample of the old population?
(Is the ‘target’ a random sample of the ‘source’?)

There are insightful comments at:

<https://mathbabe.org/2011/06/16/the-basics-of-quantitative-modeling/>

9.1 What is driving predictions? — sources of advice

The issues that arise for observational studies do not in general have clear and easy answers. The discussion on [Andrew Gelman’s blog](#)¹ canvasses some of the more important issues. There are no simple answers!

Where there are several explanatory variables, and the aim is to determine the manner in which they may be driving predictions, matters get much more complicated. Thus, in a comparison between two groups (in the example that follows, midwife led versus medical led neonatal care) one variable or factor may be of particular interest, while other variables are used to adjust for differences between the two groups that are at most a secondary focus of interest. Variables that are of secondary interest are commonly referred to as covariates. Regression coefficients can be misleading guides to what is driving predictions if one or more of the relevant covariates is not available or is not properly accounted for. A paradox of the Yule-Simpson type, sometimes referred to as Laird’s paradox, has the same potential to deceive, a potential that should be ignored.

Little that has been published since [Rosenbaum \(2002\)](#) clarifies greatly the advice that can be given for practical data analysis, beyond what Rosenbaum has to say. Note, however, that [Pearl and Mackenzie \(2018\)](#) would dispute this assessment. Pearl and his co-author do a good job of highlighting important issues that should be addressed in order to make causality judgments, at the same time overplaying what their methodology can in general achieve. If strictly implemented, the standards are so high that they severely limit what they can in practice achieve. Causality diagrams have a central role. There is a detailed, and insightful, discussion of the history that finally led to the conclusion that smoking causes lung cancer.

¹<https://statmodeling.stat.columbia.edu/2018/11/10/matching-discarding-non-matches-deal-lack-complete-overlap-regression-adjust-imbalance-treatment-control-groups/>

9.2 From air, or from water — 1849 deaths from cholera

Farr, who worked as statistician in the UK Registrar General's office, collected data on deaths from cholera in London in the 1849 epidemic. The prevailing theory at the time was that miasma, or bad air created from rotting matter, was responsible for transmitting diseases.

Farr classified districts into three groups this, according to the source of the water for most of the householders:

- 1) Thames between Battersea Bridge and Waterloo Bridge, coded as **Battersea**;
- 2) New River/Rivers Lea and Ravensbourne (sources away from the Thames), coded as **NewRiver**;
- 3) Thames between Kew and Hammersmith, i.e., further up the Thames than the first group, where the water was less polluted by sewage, coded as **Kew**.

A regression analysis, using Farr's data, gives results that have been summarized in Figure 9.1

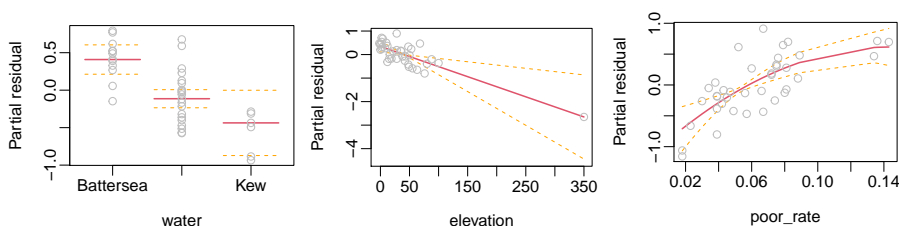


Figure 9.1: Each panel shows, in turn, the estimated contribution of a term in the model relative to the mean contribution from other model terms. Changes in deaths are on a 'log' scale, so that an increase by one unit multiplies the odds of death by close to 2.7, around an overall mean of just over six per 1000.

What can be concluded about the manner in which the three terms contributed to death rates. None of the terms stands out as being substantially more important than any other. Higher rates for the poor, who were more likely to be living in crowded conditions where it was difficult to maintain hygiene, were to be expected.

[Snow \(1855\)](#) argued that those living close to the Thames, and especially in the South, were more likely to be getting their water from or via sources that were

likely to be contaminated with human excreta. The piping of water up to higher ground gave contaminants more time to settle, with less chance of exposure to human excreta. He gave examples that he had observed directly, where the likely means of transmission of the infection appeared to be via a water source, or from poor hygiene.

Farr gave Snow's arguments some credibility, but discussed ways that the air might be the main source of transmission of an organism responsible for the disease, which multiplied in a process akin to fermentation that was presumed to take place in putrefying matter.

A context has to be provided in which to interpret the data and the regression results. While Snow had a better understanding of the contextual information, it was not comprehensive enough to persuade other medical specialists. Data from the 1854 epidemic, where it was possible to compare deaths supplied from a company that continued to get its supply from lower highly polluted Thames water with that from the company that had moved its supply higher up to less polluted water, seems in retrospect to clinch the issue. The perspective brought by germ theory would come later, with the work of Pasteur in the late 1850s and Koch in the 1880s.

9.3 Are there missing covariates?

The (Wernham et al., 2016) study used data from 244,047 singleton term deliveries that occurred between 2008 and 2012. It made the claim that midwife led care, as opposed to medical led care, gave a greater risk of adverse fetal and neonatal outcomes. Notably, the claim was that midwife led care resulted in a lower Apgar score (a measure of infant health immediately after birth) and a greater risk of the imprecisely defined diagnosis of birth asphyxia. Studies that are similarly relatively carefully done, but naive in the weight placed on the regression results, are embarrassingly common.

This study was then the basis for exaggerated claims in an article in the October 8-14 2016 issue of the NZ Listener (Chisholm, 2016, "Birth Control"). Contrary to what was claimed, the research did not "lob a grenade into the historically war-torn territory of New Zealand's maternity care." Even less did its results warrant the melodramatic claims of "Alarming maternity research" and "Revolution gone wrong" that appeared on the Listener's front cover.

9.4. THE MAY 2020 LANCET PAPER THAT WAS QUICKLY WITHDRAWN79

A major issue with the analysis is that it relies on using the [NZ Deprivation Index] (<https://www.health.govt.nz/publication/nzdep2013-index-deprivation>) to adjust for socioeconomic differences. This provides a deprivation score for meshblocks, each of around 60–110 people. It estimates the relative socioeconomic deprivation of an area, and does not directly relate to individuals. Deprived areas will often include some individuals with high socioeconomic status. Caesarean section, as a delivery type, may well have been more accessible for those of higher socioeconomic status. For National Women’s in Auckland, the elective Caesarean rate at term over 2006–2015 for doctor-led care was 32.8%, as against 7.4% for self employed midwives [farquhar2016letter]. Effects from fetal alcohol syndrome were not accounted for, nor were direct effects from substance abuse. According to NZ Ministry of Health information, international data indicates that [fetal alcohol syndrome may affect as many as 3% of births] (<https://www.health.govt.nz/our-work/diseases-and-conditions/fetal-alcohol-spectrum-disorder>)

There are analysis tools, and associated graphs, that the authors of the study could and should have used to shed light on the likely effectiveness of the covariate adjustments.

9.4 The May 2020 Lancet paper that was quickly withdrawn

Thirteen days after it was published on May 20 2020, three of the four authors withdrew a paper that claimed to find that malaria drugs, when used experimentally with patients with Covid-19, led to around 30% excess deaths. Irrespective of the problems with the data that will be noted shortly, serious flaws in the analysis ought to have attracted the attention of referees. There was inadequate adjustment for known and measured confounders (disease severity, temporal effects, site effects, dose used).

The study claimed to be based on data from 96,032 hospitalized COVID-19 patients from six continents, of which 66% were from North America. Very soon after it appeared, the article attracted critical attention, with a number of critics joining together to submit the letter [Watson et al. \(2020\)](#) to Lancet.

The sources from which the data had been obtained could not be verified, data that claimed to be from just five Australian sources had more cases than the total

of Australian government figures, and similarly for Australian deaths, there were implausibly small reported variances in baseline variables, mean daily doses of hydroxychloroquine that were 100 mg higher than US FDA recommendations.

Randomized trials designed to test the effectiveness of the drugs, and that were in progress at the time when the paper appeared, were temporarily halted. The eventual conclusion was that the drugs did not improve medical outcomes. There was some evidence that hydroxychloroquine could have adverse effects.

With current web-based technology, RCTs can be planned and carried out and yield definitive answers, in much the same time as it would take to collect and analyze the data that are required for an observational study whose conclusions can be, at best, suggestive. Data confidentiality issues are easier to handle in the context of an RCT.

9.5 The uses and traps of “algorithmic” methods — trees

Take as an example spam prediction, using tree-based methods. The boxplots in the figure show the distributions of variable values.

```
cap16 <- paste("Boxplots, showing distribution of variable values
               in data used to predict email spam")
```

```
par(oma=c(0,3,0,0), mfrow=c(2,6), mgp=c(1,0.5,0), mar=c(3.1,3.6,1.6,.6), las=0)
nam <- c("crl.tot", "dollar", "bang", "money", "n000", "make")
nr <- sample(1:dim(DAAG::spam7)[1], 500)
yesno <- DAAG::spam7$yesno[nr]
spam7 <- DAAG::spam7[nr, nam]
nam2 <- names(spam7)
nam2[1] <- "Total runs of capitals"
nam2[2] <- "No. of '$' as % of symbols"
nam2[3] <- "No. '!' as % of symbols"
nam2[4] <- "No. 'money', as % of words"
nam2[5] <- "No. 'make', as % of words"
nam2[6] <- "No. '000', as % of words"
spam7.2 <- spam7
spam7.2[,1] <- log(spam7.2[,1] + 0.5)
```


9.5. THE USES AND TRAPS OF “ALGORITHMIC” METHODS – TREES81

```

spam7.2[,2:6]<-log(spam7.2[,2:6]+0.5)
for (namtxt in nam){
  boxplot(split(spam7[,namtxt],yesno),cex=0.65,axes=F,boxwex=0.5)
  box()
  par(mgp=c(1,.5,0))
  axis(2, cex.axis=1)
  par(mgp=c(1,.25,0))
  axis(1,at=1:2,labels=c("n","y"))
  i <- match(namtxt,nam)
  mtext(side=2,line=1.75,nam2[i],adj=0.5,cex=0.8)
}
xval <-c(.1,.2,.5,1,2,5,10,20,50,100,200,500,1000,2000)
for (namtxt in nam){
  boxplot(split(spam7.2[,namtxt],yesno),cex=0.65,axes=F,boxwex=0.5)
  box()
  ranx <- range(spam7[,namtxt])
  yloc<-xval[xval>=min(ranx)&xval<max(ranx)]
  par(mgp=c(1,.5,0))
  axis(2,at=log(yloc+0.5),labels=paste(round(yloc,1)), cex.axis=1)
  par(mgp=c(1,.25,0))
  axis(1,at=1:2,labels=c("n","y"))
  i <- match(namtxt,nam)
  if(i==1)mtext(side=2,line=1,"(Logarithmic scales)",outer=T,at=0.25)
  mtext(side=2,line=1.75,nam2[i],adj=0.5,cex=0.8)
}

```

The decision tree that is obtained is:

```

cap17 <- paste("Decision tree for spam data. If the condition is satisfied, take
               the branch to the left. Otherwise, take the branch to the right.")

```

```

par(mar=c(4.1,3.1,2.6,0.6), xpd=TRUE)
require(rpart)
spam.rpart <- rpart(formula = yesno ~ crl.tot + dollar + bang +
  money + n000 + make, data=DAAG::spam7)
plot(spam.rpart, uniform=TRUE)
text(spam.rpart)
par(mar=c(0,4,0,0))
plot(c(1,8),c(0.5,8.5), asp=0.4, axes=FALSE, type="n", bty="n", cex=0.8, xlab="",yl

```

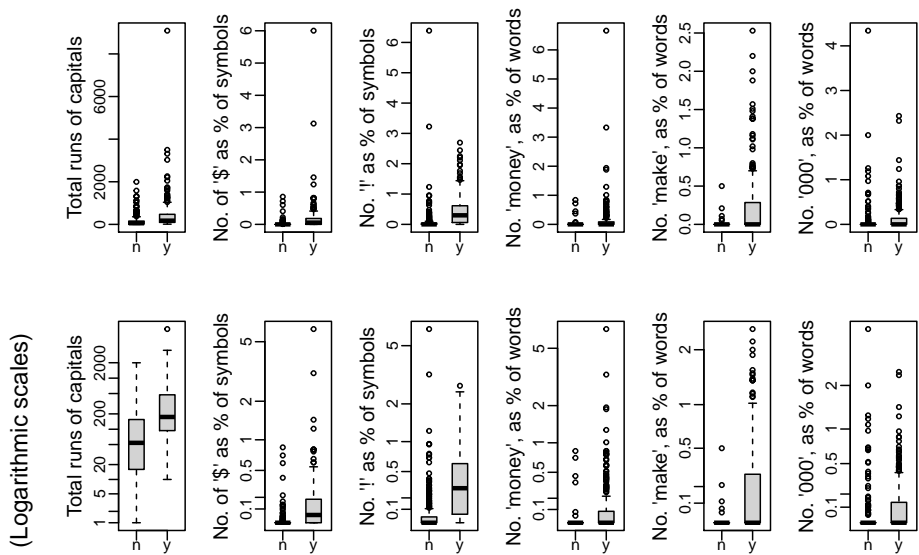
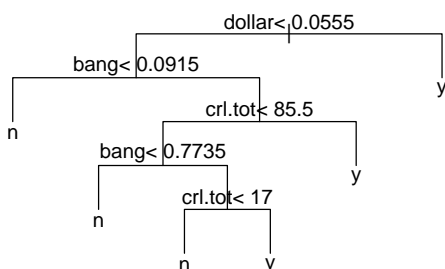


Figure 9.2: Boxplots, showing distribution of variable values in data used to predict email spam

```

text(1,8.0, " ",pos=4)
text(2,8.0, "Symbols used in tree are:",pos=4)
text(2,6.0, "dollar: Number of `$` symbols\n(as % of symbols)",pos=4)
text(2,3.5, "bang: Number of `!` symbols\n(as % of symbols)",pos=4)
text(2,1.5, "crl.tot: Total runs of capitals",pos=4)
text(2,0.75, " ", pos=4)

```



Symbols used in tree are:
dollar: Number of ‘\$’ symbols
(as % of symbols)
bang: Number of ‘!’ symbols
(as % of symbols)
crl.tot: Total runs of capitals

Figure 9.3: Decision tree for spam data. If the condition is satisfied, take the branch to the left. Otherwise, take the branch to the right.

From trees to forests

Trees such as shown will often have poor predictive power. A much more effective way to use “trees”, in many or most cases, is to make a forest (a “random forest”), and then vote between the trees. A downside is that “Random forests” and similar methods operate largely as black boxes.

- Random forest type methods may work well when the way that explanatory factors conspire to give an output is unclear.
- What works, but one does not know why, may be effective for present circumstances.
- This can be both a trap and a virtue. Thus, for detecting spam:
 - When it fails, we will likely have few clues why!
 - This may, for a short time, impede spammers!
- Spammers are anyway continually refining their strategies
 - Spam detectors must be responsive to new challenges
- Automated systems that can be easily gamed abound. They are a menace!

It helps to know the how and why of the algorithms used

Cathy O’Neill: “... it’s not enough to just know how to run a black box algorithm. You actually need to know how and why it works, so that when it doesn’t work, you can adjust.”

9.6 Regression bloopers – examples of other traps

Herrricanes vs Himmicanes

```
cap18 <- "Deaths versus damage estimate in US dollars, with logarithmic scales
          on both axes. Separate fitted lines for male and female
          hurricanes cannot be distinguished. Jung et al used a
          logarithmic scale on the vertical axis only, which on
          this graph leads to the dashed curves."
```

Authors of a paper titled “Female hurricanes are deadlier than male hurricanes” (Jung et al., 2014) compared death rates from hurricanes with female names with death rates for those given male names (jokingly called himmicanes), for 94 Atlantic hurricanes that made landfall in the United States during 1950–2012. The suggestion was that authorities took the risk from hurricanes with female names less seriously. A storm of reaction on the blogosphere was mostly himmricane!

As the primary measure of the risk posed by the hurricanes, the authors used a 2013 US\$ estimate of damage that could have been expected from a comparable hurricane in 2013. Figure 9.4 uses, instead, the estimate of damage at the time, converted to 2014 US\$.

The Jung et al analysis was roughly equivalent to regressing $\log(\text{deaths})$ on a their 2013 US\$ measure of damage, accounting also for femaleness of the name, a minor effect from barometric pressure at landfall, and interactions. Why did the authors not use, at least as a starting point, the same transformation on both axes, as in Figure 9.4? The femaleness effect then vanishes.

Jung et al’s damage measure, because designed for use for 2013 insurance purposes, assessed risk to the infrastructure that was in place in 2013. The damage

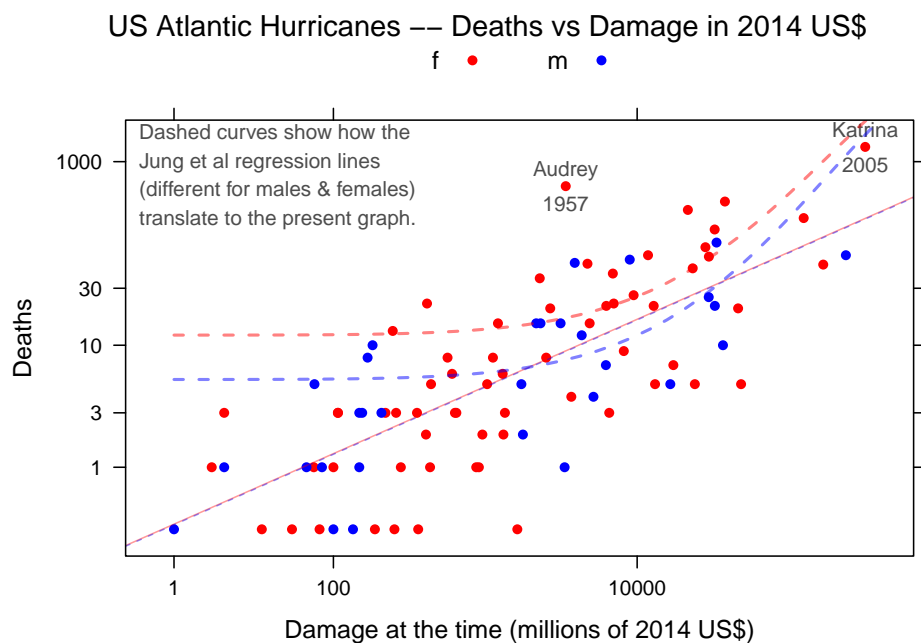


Figure 9.4: Deaths versus damage estimate in US dollars, with logarithmic scales on both axes. Separate fitted lines for male and female hurricanes cannot be distinguished. Jung et al used a logarithmic scale on the vertical axis only, which on this graph leads to the dashed curves.

measure that is more relevant to risk to human life at the time is damage caused at the time, in inflation-adjusted dollars. This makes little difference to the graph or to the analysis.

Yet another issue was that the judgment on how female a name sounded was made by students in 2014. Unconscious judgments that might have influenced disaster management, at the time when the hurricanes occurred, would very likely have been different.

Historical speed of light estimates — is there a pattern?

Creationist Barry Setterfield has argued that a reduction over time in the speed of light has led the passage of time to slow down, relative to the remote past, so that the universe is thousands rather than billions of years old. His arguments rely on making various adjustments to figures obtained historically, selecting what he regarded as the most reliable data, and then fitting a curve. He tells a story that is very different from that of Panel A of Figure 9.5. Data are from https://en.wikipedia.org/wiki/Speed_of_light.

```
cvalues <- data.frame(
  Year = c(1675, 1729, 1849, 1862, 1907, 1926, 1950, 1958, 1972),
  speed = c(220000, 301000, 315000, 298000, 299710, 299796,
            299792.5, 299792.50, 299792.4562)/1000,
  error = c(NA, NA, NA, 500, 30, 4, 3, 0.1, 0.00111)/1000
)
```

```
cap19 <- "Successive speed of light estimates.
Error estimates are available for the 1855 and later
measurments. Panel B limits attention to measurements
made in 1926 and later. The line
was fitted with no adjustment for the very different error
estimates. The dashed curve, which incorporates
such adjustments, is statistically indistinguishable
from the red horizontal line."
```

```
par(mfrow=c(1,2), mar=c(3.1,4.1,2.1,0.6))
plot(speed ~ Year, data=cvalues, cex=1.0, cex.lab=1.0, pch=1,
      xlab="", ylab="Speed (1000s of km/s)")
rect(1915,296.5,1980, 303, col="lightgray", border=NA)
```

```

with(cvalues, points(speed ~ Year, pch=16, cex=1.0))
obj <- lm(speed ~ Year, data=cvalues)
abline(obj)
mtext("A: All measurments", side=3, line=0.75, cex=1.25, at=1650, adj=0)
subdata <- subset(cvalues, Year>=1926)
ylim <- with(subdata, range(c(speed-error, speed+error), na.rm=TRUE))
plot(speed ~ Year, data=subdata, ylim=ylim, pch=0, cex.lab=1.15,
      xlab="", ylab="")
obj <- lm(speed ~ Year, data=subdata)
abline(obj)
obj3 <- lm(speed ~ poly(Year,2), data=subdata, weights=error^-2)
obj4 <- lm(speed ~ 1, data=subdata, weights=error^-2)
lines(subdata$Year, fitted(obj3), lty=2)
lines(subdata$Year, fitted(obj4), col='red')
mtext("B: From 1926, with confidence bands", side=3, line=0.75, at=1922, adj=0, cex=
with(subdata, segments(Year, speed - error, Year, speed +
  error))
with(subdata, segments(Year - 1.25, speed - error, Year +
  1.25, speed - error))
with(subdata, segments(Year - 1.25, speed + error, Year +
  1.25, speed + error))

```

Even if one were to accept Setterfield's manipulation of the data, it makes no sense at all to fit either lines such as are shown, or curves, to data values which have such very different accuracies as those shown in the graphs. For the measurements from 1862 onwards, estimates of accuracy are available. Until 1950, each new estimate lay outside the bounds for the previous estimate, indicating that these were underestimates.

The right panel is limited to the points from 1926 and on, marked off with the gray background on the left panel.

9.7 Global mean temperature trends

Figure 9.6 plots global air and sea surface temperature anomaly data against year. Anomalies, in hundredths of a degree centigrade, are differences from the 1951-1980 global average. The grey curve plots the average anomaly up to that

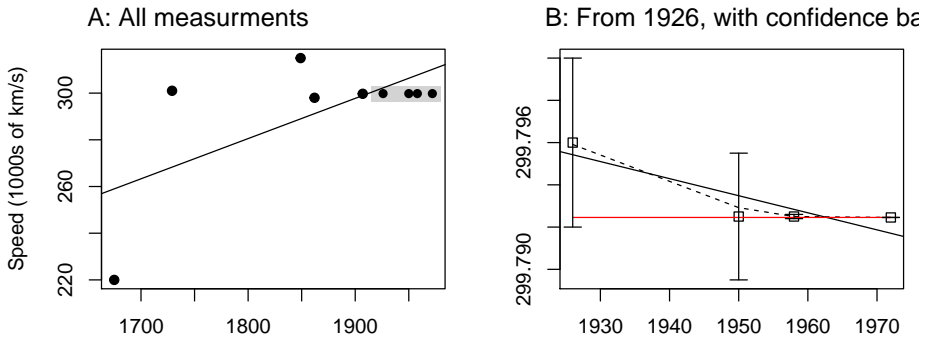


Figure 9.5: Successive speed of light estimates. Error estimates are available for the 1855 and later measurements. Panel B limits attention to measurements made in 1926 and later. The line was fitted with no adjustment for the very different error estimates. The dashed curve, which incorporates such adjustments, is statistically indistinguishable from the red horizontal line.

point in time.

cap20 <- "Anomalies (differences) in hundredths of a degree centigrade from global average temperatures over 1951-1980, plotted against year. The gray curve shows, for each year, the average anomaly up to that point in time. The last year in which this lay below the gray line was 1962."

```
## ---- loti-smooth -----
load('data/loti.RData')
anomaly <- loti[, "J.D"]
num <- seq(along = anomaly)
AVtodate <- cumsum(anomaly)/num
yr <- loti$Year
anomTxt <- "Difference from baseline"
yl = substitute(txt * " (0.01" * degree * "C)",
                list(txt=anomTxt))
plot(yr, anomaly, xlab = "Year", ylab = as.expression(yl))
mtext(side=3, line=0.75,
      "Global temperature differences from 1951-1980 global average")
```



```

lines(AVtodate ~ yr, col = "gray", lwd = 2)
lastLessYr <- max(yr[anomaly < AVtodate])
lastLessYr <- loti[as.character(lastLessYr), "J.D"]
yarrow <- lastLessYr - c(4, 0.75) * strheight("0")
arrows(lastLessYr, yarrow[1], lastLessYr, yarrow[2],
       col = "gray", lwd = 2)

```

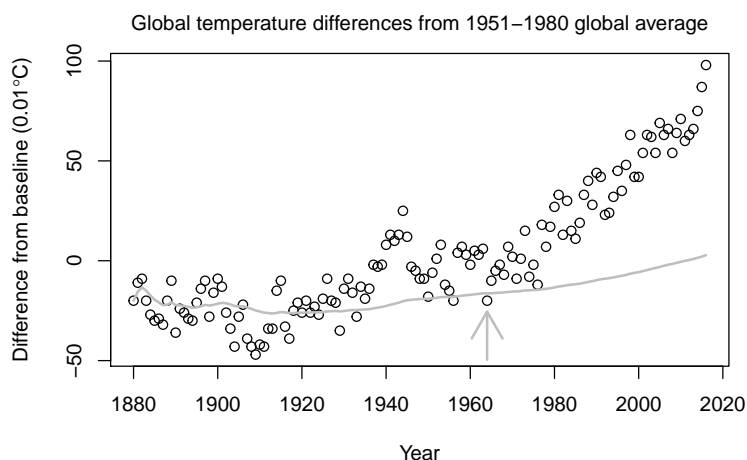


Figure 9.6: Anomalies (differences) in hundredths of a degree centigrade from global average temperatures over 1951–1980, plotted against year. The gray curve shows, for each year, the average anomaly up to that point in time. The last year in which this lay below the gray line was 1962.

Observe that 1964 was the last year in which the global temperature fell below the average to that time. For the 52 subsequent years (from 1965 to 2016 inclusive), the global average was above the average up to that date. Under the (false) assumption that global temperature is varying randomly (and therefore independently) about a common mean, the probability of this happening is $2^{-40} = 9.1 \times 10^{-13}$. A variation of this argument came from a speaker on the Australian ABC Science Show on April 3 2011. Under any model that accounts for what are now fairly well understood patterns of correlation over time, the probability, while very small, is not that small! Arguments that overstate the case for what is now a well-established pattern of change are unhelpful

It is likewise nonsensical to fit a line to the cherry-picked years 1998-2008, where the trend was relatively flat.

Chapter 10

The critique of scientific claims

An ideal is that scientific processes will keep to a minimum the publication of claims that are not supported by evidence. Like anything designed by humans, and put into practice by humans, this is not precisely how it works in practice. Notwithstanding peer review processes that are designed to ensure that claims made in scientific papers withstand careful scrutiny, the traps and fallacies described in these notes, and more besides, are found in the scientific literature. The effectiveness of the peer review process varies widely from one area of science to another.

10.1 What results can be trusted?

Scientific processes work best when claims made by one scientist or group of scientists attract widespread interest and critique from the wider group of scientists who work in the same general area. Examples are the May 2020 Lancet and New England Journal of Medicine papers arguing that use of the drug hydroxychloroquine as a treatment for Covid-19 was increasing patient deaths. Issues with these papers were quickly identified because they made claims that bore on an issue of major concern, and attracted attention from readers who

carefully scrutinized their detailed statements. They were quickly retracted.

Heavy reliance on the sharing of data and skills, and full use of the benefits that modern technology has to offer, have been vital to progress in such areas as earthquake science, the study of viruses and vaccines, modelling of epidemics, and climate science. This sharing of data and skills, and use of modern technology, also helps in the critique of what has been published earlier. Areas where there has not been the same impetus for change are much more susceptible to the damage that arises from systems for funding and publishing science that encourage the formal publication of what would better be treated as preliminary results — a first stab at an answer. Publication of experimental results should not be a once-for-all event, but a staged process that moves from “this looks promising” to “has been independently replicated”, and to post=publication critique.

Publication does not of itself validate scientific claims, Rather, as stated in [Popper \(1963\)](#)

Observations or experiments can be accepted as supporting a theory (or a hypothesis, or a scientific assertion) only if these observations or experiments are severe tests of the theory – or in other words, only if they result from serious attempts to refute the theory.

10.2 A look at what can go wrong

Fraud, though uncommon, happens more often than one might hope. What is disturbing is the small number of scientists with large numbers of papers that were retracted on account of fraud. How were they able to get away with publishing so many papers, usually with fraudulent data, before the first identification of fraud that led to a checking of all their work? [Ritchie \(2020\)](#) (pp67-68) cites, as an extreme example, the case of a Japanese anesthesiologist with 183 retracted papers.

More common are mistakes in data collection, unacknowledged sources of bias, hype, mistakes or biases in the handling of data and/or in data analysis, attaching a much higher degree of certainty to statistical evidence than the results warrant, and selection effects. Selection effects can work in various ways — selection of a subset of data where there appears to be an effect of interest, choice

of the outcome variable and/or analysis approach that most nearly gives the result that is wanted, and so on. In analysis of data from experiments where two treatments are compared, the common use of the arbitrary $p \leq 0.05$ criterion (see the next subsection) as a cutoff for deciding what will be published has the inevitable effect of selecting out, in contexts where there was no difference of consequence one in twenty of such results for publication. This adds to other selection effects.

10.3 The case of Eysenck and his collaborators

At the time of his death in 1997, Eysenck was the living psychologist most frequently cited in the peer-reviewed scientific literature. Much of his work was controversial in its time, with papers containing “questionable data and results so dramatic they beggared belief” [o2020famous]. He relied heavily on what has now been identified heavily doctored data that was supplied to him by German collaborator Grossarth-Maticek. Particularly egregious was the claim that individuals with an identifiably cancer-prone personality had a risk of dying from cancer that was as much as 121 times higher than that of people with a “healthy” personality — one of several links that the duo claimed to have found between personality and mortality. Investigations into Eysenck’s work, including collaborative work with Grossarth-Maticek, are ongoing. Fourteen papers have been retracted, and another 71 have received “expressions of concern”. A large replication study conducted in 2004 found none of the claimed links, apart from a modest link between personality and cardiovascular disease.

10.4 The use of artificial intelligence to detect Covid-19

[Roberts et al. \(2021\)](#) identified an astonishing 320 papers and preprints that appeared between 1 January 2020 to 3 October 2020, and which describe the use of new machine learning models for the diagnosis of COVID-19 from chest radiographic (CXR) and chest computed tomography (CT) images. Quality screening reduced this number to 62, which were then examined in more detail. None of the 62 satisfied these more detailed requirements, designed to check whether the algorithms used had been shown to be effective for use in clinical

practice. Among other deficiencies, 48 did not complete any external validation, and 55 had a high risk of bias with respect to at least one of participants, predictors, outcomes and analysis. In a popular account of the results that appeared in *New Scientist*, [Roberts \(2021\)](#) comments that, relative to persisting to develop a model that will survive a rigorous checking process and might be used in practice, “it is far easier to develop a model with poor rigour and [apparent] excellent performance and publish this.” This is a damning indictment of the way that large parts of the research and publication process currently work. The public good would be much better served by a process that encourages researchers to persist until it has been demonstrated that researchers have a model that meets standards such as are set out in [Roberts et al. \(2021\)](#). One may hope that one result of this work will be to shift the research and publication focus accordingly.

10.5 Laboratory experimental science — what do we find?

In the past several years, there has been a steady accumulation of evidence that relates to the claim, in [Ioannidis \(2005\)](#), that “most published research findings are false”. Ioannidis has in mind, not published results in general, but primarily laboratory studies. Papers that have added to the body of evidence that broadly support claims made in the Ioannidis paper include:

- Amgen: Reproduced 6 only of 53 ‘landmark’ cancer studies.
 - [Begley and Ellis \(2012\)](#)
 - [Begley \(2013\)](#) notes issues with the studies that failed
- Bayer: Main results from 19 of 65 ‘seminal’ drug studies
 - NB, journal impact factor was not a good predictor!
 - [Prinz et al. \(2011\)](#)
- fMRI studies: 57 of 134 papers (42%) had ≥ 1 case lacking check on separate test image. Another 14%, unclear ...
 - [Kriegeskorte et al. \(2009\)](#)

The psychological science community is further advanced in addressing these issues that many other communities, with The Center for Open Science (COS) taking a strong lead in studies designed to document the extent of the issues.

Other Center for Open Science (COS) Projects have been:

- Many Labs — reproduce 13 classical psych studies
 - Of 13 studies — 10: successful, 1: weakly, 2: no!
 - Plots show scatter across the 36 participating teams
 - [Klein et al. \(2014\)](#)
- Cancer Studies — 50 “most impactful” from 2010-2012
 - [Kaiser \(2015\)](#)

Details from other relevant studies are given in the recent book [Ritchie \(2020\)](#) “Science fictions: Exposing fraud, bias, negligence and hype in science.” Problems arise, primarily, in areas where relatively small groups of scientists, with similar training and skills, work independently. The critiques have limited relevance to areas where the nature of the work forces collaboration between scientists with diverse skills, widely across different research groups.

Even with such checks as there are on published research, it is clear that not everything that passes the peer review process can be relied on. The case is far worse for claims made, often for reasons of commercial gain, who have no evidence at all for the effectiveness of the therapy or remedy that they promote.

10.6 The Reproducibility: Psychology project

Figure [10.1](#) summarizes evidence from the Reproducibility: Psychology project ([OSC, 2015](#)). The effect size is the difference between the two means that are to be compared, divided by the pooled standard deviation for the two groups.

- It chose 100 studies (3 journals, 2008)
- There was one replicate only of each study
- Initial results: 39 agree + 24 ? “similar” + 37 fail

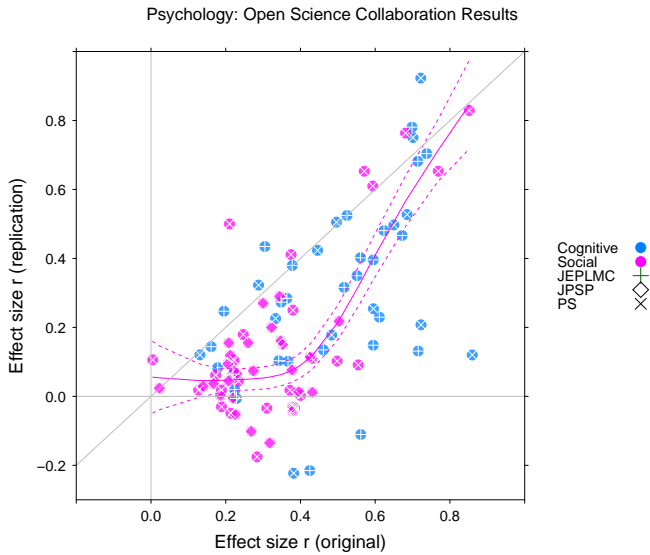


Figure 10.1: Psychology reproducibility project — Effect sizes are compared between the replication and the initial study.

Notice that the effect size is almost always smaller for the replication. A smooth curve, with confidence interval, has been fitted for results in social psychology. Notice that, for social psychology, it is only for an original effect size greater than 0.4 that one starts to see a positive correlation between the effect sizes for replicate and original.

10.7 When science and commercial interests collide

The tobacco and fossil fuel industries

The tobacco industry has made extensive use of “the science is not settled” arguments in its efforts to dismiss the evidence that smoking causes lung cancer. The goal has been to raise doubt, create confusion and undermine the science.

The same PR firms, and the same researchers used to support these efforts were later used in the attempt to undermine climate science.¹

In spite of the support that has been available from the fossil fuel industry for research that is critical of mainstream climate science research, no substantially different alternative account has emerged, and no climate change models have emerged that give results that are widely different from the consensus. An interesting case is that of Richard Muller's Berkeley Earth Surface Temperature Study (BEST), with a large part of the funding coming from the right-wing billionaire Charles Koch, known for funding climate skeptic groups such as the Heartland Institute.

Richard Muller had been known for his skepticism, in part supported by legitimate scientific concerns. He made headlines when he announced his acceptance of what climate scientists had already been saying for more than 15 years previous. > After years of denying global warming, physicist Richard Muller now says "global warming is real and humans are almost entirely the cause."²

The broad sweep of work in climate science is, because it has survived informed critique, and because of the diversity of contributing skills and data, unusually secure. Details, especially as they affect what may happen in individual countries, are subject to continual revision.

A particular issue is the role of the greenhouse gases, and of their interactions with water vapour. Much larger than their direct effect is the contribution that comes from their warming of the air in which water vapour is present, allowing it to retain more vapour and trap more heat). A standard denialist trump card has been to claim the authority of scientists who have standing in their own areas for the claim that the contribution of greenhouse gases is inconsequential relative to that of water vapour.³ The scientists involved are among a number who have fossil fuel industry links, and have been willing to allow themselves to be used as advocates for an industry that sees action on climate change as a threat.⁴

¹<https://www.scientificamerican.com/article/tobacco-and-oil-industries-used-same-researchers-to-sway-public1/>>

²<https://courses.seas.harvard.edu/climate/eli/Courses/global-change-debates/Sources/Hockeystick-global-temperature/more/Richard-Muller/Muller-is-a-believer-Hallelujah.pdf>>

³See <https://west.web.unc.edu/climate-change/> and <>

⁴<https://insideclimatenews.org/news/12032015/leaked-email-reveals-whos-who-list-climate-denialists-merchants-of-doubt-oreskes-fred-singer-marc-morano-steve-milloy>

Big Pharma

Between 1999 and 2019, opioid deaths in the United States increased by a factor of six, to almost 50,000.

A large contributor has been the increased use of prescription opioids. Purdue Pharma stands out for its aggressive marketing of oxycodone, sold under the brand name OxyContin, arguing that concerns over addiction and other dangers from the drugs were overblown. In September 2019, Purdue Pharma declared bankruptcy, facing significant liability in OxyContin and opioid addiction lawsuits. Details of settlements are still being worked in the courts.⁵

Strict regulations govern, in most countries, the approval of prescription drugs. Purdue Pharma exploited the much more limited control over off-label use of approved drugs, i.e., use for purposes for which they have not received formal approval, and for a time stayed under the radar.

Another scandal that demonstrates how drug companies can sometimes work their way around the approval process concerns the drug Vioxx, likewise marketed as a painkiller.(Valentine and Prakash, 2007) Concerns that the drug might be increasing the risk of heart attacks began to emerge in the months following its approval for use in May 1999. By November 1, a study set up to investigate these concerns reported 79 heart attacks out of 4000 among those taking the drug, as opposed to 41 among a comparator group that was taking naproxen. The drug continued on the market as the evidence against Vioxx strengthened further. It was argued, with no evidence to back this up, that naproxen likely had a protective effect. In any case, why allow Vioxx to go to market when naproxen was clearly carried less risk.

In September 2004, when a colon-polyp prevention study showed that Vioxx increased the risk of heart attack after 18 months, Merck withdrew the drug. A Lancet paper that was published later estimated that between 88,000 and 140,000 Americans had heart attacks from taking Vioxx.⁶ The increased risk continued long after patients had ceased taking the drug.

⁵ <https://topclassactions.com/lawsuit-settlements/open-lawsuit-settlements/opioids/purdue-opioid-addiction-class-action-settlement/>

⁶Jüni et al. (2004)

What can one trust?

One may hope that checks on any new drug will be stricter than was the case for Vioxx, that lessons have been learned, that where in future drugs come up for approval, checks will continue on possible long-term effects. With fringe and quack medicines, there are no such checks.

It has been interesting to follow approval processes for Covid-19 vaccines. At the time of writing, the Pfizer, Astrazena and Moderna vaccines have had, in addition to their testing in clinical trials, extra-ordinary levels of testing in clinical practice, with risks that are small relative to the small risk that we take ever time we cross a busy road.

10.8 Tricks used to dismiss established scientific results

The headings, and some of commentary, are adapted from an article by Associate Professor Hassan Vally from La Trobe University, [that appeared in *The Conversation*](https://theconversation.com/5-ways-to-spot-if-someone-is-trying-to-mislead-you-when-it-comes-to-science-138814?utm_medium=email&utm_campaign=Latest%20from%20The%20Conversation%20for%20March%209%202021%20-%201883318379&utm_content=Latest%20from%20The%20Conversation%20for%20March%209%202021%20-%201883318379+CID_6009c8d7af2a01f376a88d0491598cfa&utm_source=campaign_monitor&utm_term=5%20ways%20to%20spot%20if%20someone%20is%20trying%20to%20mislead%20you%20when%20it%20comes%20to%20science)⁷

1. “The ‘us versus them’ narrative”
The powers that be are trying to deceive us.” Ask who is making real attempt to deceive — commercial interest groups, peddlers of quack treatments, influence peddlers, . . .
2. ‘I’m not a scientist, but...’ Meaning perhaps: “I’m not a scientist, but that does not stop me making an authoritative pronouncement that flies in the face of established results.” Or: “I know what the science says, but I’m keeping an open mind”.
Politicians are among the most frequent offenders.
3. Reference to ‘the science not being settled’
There are of course times when the science is not settled, and when this

⁷https://theconversation.com/5-ways-to-spot-if-someone-is-trying-to-mislead-you-when-it-comes-to-science-138814?utm_medium=email&utm_campaign=Latest%20from%20The%20Conversation%20for%20March%209%202021%20-%201883318379&utm_content=Latest%20from%20The%20Conversation%20for%20March%209%202021%20-%201883318379+CID_6009c8d7af2a01f376a88d0491598cfa&utm_source=campaign_monitor&utm_term=5%20ways%20to%20spot%20if%20someone%20is%20trying%20to%20mislead%20you%20when%20it%20comes%20to%20science

is the case, one can expect scientists to openly argue different points of view based on the evidence available. Or, what has come to be accepted wisdom may turn out to be wrong or in need of substantial revision. But challenges to the accepted wisdom have to be carefully argued, and themselves survive informed critique.

4. Overly simplistic explanations

Oversimplifications and generalizations are central to many conspiracy theory arguments. Science is often messy, complex and full of nuance. The truth can be much harder to explain, and can sometimes sound less plausible, than a simple but incorrect explanation. Use of simplistic statistical arguments is common, e.g., fit a line to the cherry-picked years 1998-2008, ignoring year to year correlation as well as influences that operate over longer time periods.

5. Cherry-picking

One is not entitled to choose one study over another just because it aligns with what you prefer to believe. This is not how science works.

Not all studies are equal; some provide much stronger evidence than others. The way that choices are made has to stand up to critical scrutiny.

Further reading — books, videos, and websites

The books noted have all been referred to in the text.

- Kahneman (2013) . Thinking, fast and slow.
 - [Interview with Kahneman](#)⁸
 - [A brief animated overview of some key points](#)⁹
- Smith (2014) . Standard Deviations: Flawed Assumptions, Tortured Data, and Other Ways to Lie with Statistics
 - [Brian Lehrer interview with Smith](#)¹⁰
- Ellenberg (2015) . How not to be wrong.
 - [There are links to several Ellenberg video clips](#)¹¹
- Levitin (2016) . A field guide to lies and statistics.
- Nisbett (2016) . Tools for smart thinking.
- Cairo (2013) . The functional art: an introduction to information graphics and visualization.
- Ritchie (2020) . Science fictions: Exposing fraud, bias, negligence and hype in science.
- [BBC links to helpful web resources](#)¹²

⁸<https://www.youtube.com/watch?v=PirFrDVRBo4>

⁹<https://www.youtube.com/watch?v=uqXVAo7dVRU&t=28s>

¹⁰<http://www.garysmithn.com/standard-deviations.html>

¹¹<http://www.thelavinagency.com/news/new-videos-jordan-ellenberg-runs-the-numbers>

¹²<http://www.bbc.co.uk/editorialguidelines/guidance/reporting-statistics>

Bibliography

- Arkes, H. R. and Gaissmaier, W. (2012). Psychological research and the prostate-cancer screening controversy. *Psychological science*, 23(6):547–553.
- Begley, C. G. (2013). Reproducibility: Six red flags for suspect work. *Nature*, 497(7450):433–434.
- Begley, C. G. and Ellis, L. M. (2012). Drug development: Raise standards for preclinical cancer research. *Nature*, 483(7391):531–533.
- Boot, H. and Maindonald, J. (2008). New estimates of age-and sex-specific earnings and the male–female earnings gap in the british cotton industry, 1833–1906 1. *The Economic History Review*, 61(2):380–408.
- Brawley, O. W. (2018). Prostate cancer screening: And the pendulum swings.
- Cairo, A. (2013). *The functional art: an introduction to information graphics and vizualisation*. New Riders, 1 edition.
- Chisholm, D. (2016). Birth control. *The NZ Listener*, October 8 - 14, pages 18–24. October 8 - 14.
- Coleman, T. (2019). Causality in the time of cholera: John snow as a prototype for causal inference. *Available at SSRN 3262234*.
- Collins, J. (2001). *Good to great: why some companies make the leap ... and others dont*. Random House.
- Du Toit, G., Roberts, G., Sayre, P. H., Bahnson, H. T., Radulovic, S., Santos, A. F., Brough, H. A., Phippard, D., Basting, M., Feeney, M., et al. (2015). Randomized trial of peanut consumption in infants at risk for peanut allergy. *N Engl J Med*, 372:803–813.

- Ellenberg, J. (2015). *How not to be wrong*. Penguin Books, 1 edition.
- Hassall, A. H. (1850). Memoir on the organic analysis or microscopic examination of water: Supplied to the inhabitants of london and the suburban districts. *The Lancet*, 55(1382):230–235.
- Herndon, T., Ash, M., and Pollin, R. (2014). Does high public debt consistently stifle economic growth? a critique of reinhart and rogooff.
- Ioannidis, J. P. A. (2005). Why most published research findings are false. *CHANCE*, 18(4):40–47.
- Jung, K., Shavitt, S., Viswanathan, M., and Hilbe, J. M. (2014). Female hurricanes are deadlier than male hurricanes. *Proceedings of the National Academy of Sciences*, 111(24):8782–8787.
- Jüni, P., Nartey, L., Reichenbach, S., Sterchi, R., Dieppe, P. A., and Egger, M. (2004). Risk of cardiovascular events and rofecoxib: cumulative meta-analysis. *The lancet*, 364(9450):2021–2029.
- Kahneman, D. (2013). *Thinking, fast and slow*. Farrar, Straus and Giroux, 1 edition.
- Kaiser, J. (2015). The cancer test. *Science*, 348(6242):1411–1413.
- Klein, R. A., Ratliff, K. A., Vianello, M., Adams Jr, R. B., Bahník, Š., Bernstein, M. J., Bocian, K., Brandt, M. J., Brooks, B., Brumbaugh, C. C., et al. (2014). Investigating variation in replicability. *Social Psychology*, 45(3):142–152.
- Kriegeskorte, N., Simmons, W. K., Bellgowan, P. S. F., and Baker, C. I. (2009). Circular analysis in systems neuroscience: the dangers of double dipping. *Nature Neuroscience*, 12(5):535–540.
- Levitin, D. J. (2015). *The organized mind*. Penguin, 1 edition.
- Levitin, D. J. (2016). *A field guide to lies and statistics*. Penguin Random House, 1 edition.
- Lord, F. M. (1967). A paradox in the interpretation of group comparisons. *Psychological bulletin*, 68(5):304.
- National Science Foundation, Washington, D. N. S. B. (1975). *Science Indicators, 1974*. Superintendent of Documents.
- Nisbett, R. E. (2016). *Mindware. Tools for smart thinking*. Penguin Books, 1 edition.

- OSC (2015). Estimating the reproducibility of psychological science. *Science*, 349(6251):aac4716–aac4716. Open Science Collaboration (Nosek, B A and others).
- Pearl, J. and Mackenzie, D. (2018). *The book of why: the new science of cause and effect*. Basic Books.
- Popper, K. R. (1963). Science: Problems, aims, responsibilities. *Federation of American Societies for Experimental Biology*, 22:961–972. Reprinted in Popper, K. (1994). *The Myth of the Framework: In Defense of Science and Rationality*, 82–111. (ed. M. A. Notturmo). London and New York: Routledge.
- Prinz, F., Schlange, T., and Asadullah, K. (2011). Believe it or not: how much can we rely on published data on potential drug targets? *Nature Reviews Drug Discovery*, 10(9):712–712.
- Raichand, S., Dunn, A. G., Ong, M.-S., Bourgeois, F. T., Coiera, E., and Mandl, K. D. (2017). Conclusions in systematic reviews of mammography for breast cancer screening and associations with review design and author characteristics. *Systematic reviews*, 6(1):1–8.
- Reinhart, C. M. and Rogoff, K. S. (2010). Growth in a time of debt. *American economic review*, 100(2):573–78.
- Ritchie, S. (2020). *Science fictions: Exposing fraud, bias, negligence and hype in science*. Random House.
- Roberts, M. (2021). Machine churning.
- Roberts, M., Driggs, D., Thorpe, M., Gilbey, J., Yeung, M., Ursprung, S., Aviles-Rivero, A. I., Etmann, C., McCague, C., Beer, L., et al. (2021). Common pitfalls and recommendations for using machine learning to detect and prognosticate for covid-19 using chest radiographs and ct scans. *Nature Machine Intelligence*, 3(3):199–217.
- Rosenbaum, P. R. (2002). *Observational Studies*. Springer, 2 edition.
- Simon, H. A. (1992). What is an “explanation” of behavior? *Psychological science*, 3(3):150–161.
- Smith, G. (2014). *Standard Deviations: Flawed Assumptions, Tortured Data, and Other Ways to Lie with Statistics*. Duckworth Overlook.
- Snow, J. (1855). On the mode of communication of cholera.

- Thaler, R. H. and Ganser, L. (2015). *Misbehaving: The making of behavioral economics*. WW Norton New York.
- Valentine, V. and Prakash, S. (2007). Timeline: The rise and fall of viox. *NPR (National Public Radio)*, 10.
- Waterman, R. H. and Peters, T. J. (1982). *In search of excellence: Lessons from America's best-run companies*. New York: Harper & Row.
- Watson, J., Adler, A., Agweyu, A., et al. (2020). Open letter to mr mehra, ss desai, f ruschitzka, and an patel, authors of "hydroxychloroquine or chloroquine with or without a macrolide for treatment of covid-19: a multinational registry analysis". *Lancet*, pages 31180–6.
- Wernham, E., Gurney, J., Stanley, J., Ellison-Loschmann, L., and Sarfati, D. (2016). A comparison of midwife-led and medical-led models of care and their relationship to adverse fetal and neonatal outcomes: A retrospective cohort study in new zealand. *PLOS Medicine*, 13(9):e1002134.
- Young, S. W. (2014). Improving library user experience with a/b testing: Principles and process. *Weave: Journal of Library User Experience*, 1(1).